

The Restricted Isometry Property for Random Block Diagonal Matrices

Armin Eftekhari, Han Lun Yap, Christopher J. Rozell, and Michael B. Wakin*

October 2012

Abstract

In Compressive Sensing, the Restricted Isometry Property (RIP) ensures that robust recovery of sparse vectors is possible from noisy, undersampled measurements via computationally tractable algorithms. It is by now well-known that Gaussian (or, more generally, sub-Gaussian) random matrices satisfy the RIP under certain conditions on the number of measurements. Their use can be limited in practice, however, due to storage limitations, computational considerations, or the mismatch of such matrices with certain measurement architectures. These issues have recently motivated considerable effort towards studying the RIP for structured random matrices. In this paper, we study the RIP for block diagonal measurement matrices where each block on the main diagonal is itself a sub-Gaussian random matrix. Our main result states that such matrices can indeed satisfy the RIP but that the requisite number of measurements depends on certain properties of the basis in which the signals are sparse. In the best case, these matrices perform nearly as well as dense Gaussian random matrices, despite having many fewer nonzero entries.

Keywords— Compressive Sensing, Block Diagonal Matrices, Restricted Isometry Property

1 Introduction

Many interesting classes of signals have a low-dimensional geometric structure that can be exploited to design efficient signal acquisition and recovery methods. The emerging field of Compressive Sensing (CS) deals with signals that can be parsimoniously expressed in a basis or a dictionary. A canonical result in CS states that sparse signals, i.e., signals with very few nonzero entries, can be accurately recovered from a small number of linear measurements by solving a tractable convex optimization problem if the measurement system satisfies the Restricted Isometry Property (RIP) [4].

The RIP requires the linear measurement system to approximately maintain the distance between any pair of sparse signals in the measurement space, implying that the geometry of the

*The first and second authors contributed equally to this paper. AE and MBW are with the Department of Electrical Engineering and Computer Science at the Colorado School of Mines. HLY and CJR are with the School of Electrical and Computer Engineering at the Georgia Institute of Technology. Email: aeftekha@mines.edu and yhanlun@dso.org.sg. This work was partially supported by NSF grants CCF-0830456 and CCF-0830320, by NSF CAREER grant CCF-1149225, and by DSO National Laboratories of Singapore. A preliminary version of Theorem 1, with a different proof, was originally presented at the 2011 IEEE Conference on Information Sciences and Systems (CISS) [31].

family of sparse signals is approximately preserved in the measurement space. Apart from playing a central role in the analysis of numerous signal recovery algorithms in CS [5, 19, 11], the RIP also provides a framework to analyze signal processing and inference algorithms in the compressed measurement domain [10]. Moreover, measurement systems that satisfy the RIP, after undergoing some minor modifications, can approximately preserve the geometry of an arbitrary point cloud (as confirmed by the Johnson-Lindenstrauss lemma) [17] or a low-dimensional compact manifold [33].

A measurement system represented by a matrix populated with i.i.d. sub-Gaussian¹ random variables is known to satisfy the RIP with high probability whenever the number of rows scales linearly with the sparsity of the signal and logarithmically with the length of the signal [4]. Such matrices are also universal in that, with the same number of random measurements, they satisfy the RIP with respect to any fixed sparsity basis with high probability. We refer to such matrices—densely populated with i.i.d. random entries—as *unstructured* measurement matrices. There has been significant recent interest in studying *structured* measurement systems because unstructured random measurements may be undesirable due to memory limitations, computational costs, or specific constraints in the data acquisition architecture. Many structured systems have been studied in the CS literature, including subsampled bounded orthonormal systems [24, 21], random convolution systems (described by partial Toeplitz [14] and circulant matrices [22, 16]) and deterministic matrix constructions [12]. Generally, structured random matrices require more measurements to satisfy the RIP and lack the universality of unstructured random matrices.

In this paper, we are concerned with establishing the RIP for block diagonal matrices populated with i.i.d. sub-Gaussian random variables. The advantages of such matrices are varied. First, these matrices require less memory and computational resources than their unstructured counterparts. Second, they are particularly useful for representing acquisition systems with architectural constraints that prevent global data aggregation. For example, this type of architecture arises in distributed sensing systems where communication and environmental constraints limit the dependence of each sensor to only a subset of the data and in streaming applications where signals have data rates that necessitate operating on local signal blocks rather than on the entire signal simultaneously. In these scenarios, the data may be divided naturally into discrete blocks, with each block acquired via a local measurement operator.

To make things concrete, for some positive integers J, N , and M , set $\widetilde{M} := JM$ and $\widetilde{N} := JN$. We model a signal $x \in \mathbb{C}^{\widetilde{N}}$ as being partitioned into J blocks of length N , i.e., $x = [x_1^T, \dots, x_J^T]^T$ where $x_j \in \mathbb{C}^N$, $j \in [J]$. Here, $[J]$ denotes the set $\{1, 2, \dots, J\}$. As an example, x can be a video sequence and $\{x_j\}$, $j \in [J]$, can be the individual frames in the video. For each $j \in [J]$, we suppose that a linear operator $\Phi_j : \mathbb{C}^N \rightarrow \mathbb{C}^M$ collects the measurements $y_j = \Phi_j x_j$. In our example, this means that each video frame x_j is measured with an operator Φ_j . Concatenating all of the measurements into a vector $y \in \mathbb{C}^{\widetilde{M}}$, we then have

$$\underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_J \end{bmatrix}}_{y: \widetilde{M} \times 1} = \begin{bmatrix} \Phi_1 & & & \\ & \Phi_2 & & \\ & & \ddots & \\ & & & \Phi_J \end{bmatrix} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_J \end{bmatrix}}_{x: \widetilde{N} \times 1}, \quad (1)$$

¹Roughly speaking, the tail of a sub-Gaussian random variable is similar to that of a Gaussian random variable. This term is defined precisely in Section 2.

Thus we see that the overall measurement operator relating y to x will have a block diagonal structure. In this paper we consider the two scenarios. When the $\{\Phi_j\}$ are distinct, in which case we call the resulting matrix a Distinct Block Diagonal (DBD) matrix. When the $\{\Phi_j\}$ are all identical, in which case we call the resulting matrix a Repeated Block Diagonal (RBD) matrix. Our results show that whenever the total number of measurements \widetilde{M} is sufficiently large, DBD and RBD matrices can both satisfy the RIP. As we summarize in Sections 1.2 and 1.3 below, the requisite number of measurements depends on the type of matrix (DBD or RBD) and on the basis in which x has a sparse expansion. We also show that certain sparse matrices and random convolution systems considered in the CS literature can be studied in the framework of block diagonal matrices.

In general, proving the RIP for structured measurement systems requires analytic tools beyond the elementary approaches that suffice for unstructured matrices. For example, in [24] the authors employed tools such as Dudley's inequality from the theory of probabilities in Banach spaces, and in [22] a variant of Dudley's inequality for chaos random processes was used to obtain a result that was out of reach for elementary approaches. While some of the ideas and techniques utilized in [24] and [22] can be used to establish the RIP for random block diagonal matrices (see [31] for our preliminary study), even these sophisticated tools result in measurement rates that are worse than what we report in this paper. Fortunately, recent work by Krahmer et al. has established an improved bound on the suprema of chaos random processes that enabled them to prove the RIP for Toeplitz matrices with an optimal number of measurements [16]. The bound in [16] is very general, and we have leveraged this result for the main results of this paper. Specifically, the work in [16] has allowed us to develop a unified treatment of DBD and RBD matrices with bounds on the measurement rates that are significantly improved over our preliminary work.

1.1 Definition of the RIP

A linear measurement operator satisfies the RIP if it acts as an approximate isometry on all sufficiently sparse signals. More specifically, the Restricted Isometry Constant (RIC) of a matrix $A \in \mathbb{R}^{\widetilde{M} \times \widetilde{N}}$ is defined as the smallest positive number δ_S for which

$$(1 - \delta_S)\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + \delta_S)\|x\|_2^2 \quad \text{for all } x \text{ with } \|x\|_0 \leq S, \quad (2)$$

where $\|\cdot\|_0$ merely counts the number of nonzero entries of a vector. In many applications, however, signals may be sparse in an orthobasis U other than the canonical basis, and so we will find the notion of the U -RIP more convenient.

Definition 1. Let U denote an orthobasis for $\mathbb{C}^{\widetilde{N}}$. The RIC of a matrix $A \in \mathbb{R}^{\widetilde{M} \times \widetilde{N}}$ in the basis U , $\delta_S = \delta_S(A, U)$, is defined as the smallest positive number for which

$$(1 - \delta_S)\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + \delta_S)\|x\|_2^2 \quad \text{for all } x \text{ with } \|U^*x\|_0 \leq S, \quad (3)$$

where U^* denotes the conjugate transpose of U .

In general, whenever the sparsity basis is clear from the context, U is dropped from our notation (after its first appearance). More generally, the notion of the RIP could be extended to the class of signals that are sparse in an overcomplete dictionary [6]. Nonetheless, for the clarity of exposition, we restrict ourselves throughout this paper to considering only orthobases.

1.2 The RIP for Distinct Block Diagonal (DBD) Matrices

Suppose the matrices $\{\Phi_j\}$, $j \in [J]$, in (1) are distinct, and let Ψ denote the resulting block diagonal matrix. Following [20], we say that Ψ has a DBD structure. DBD matrices arise naturally when modeling the information captured by individual (and different) sensors in a sensor network. In this setting, $y_j = \Phi_j x_j$ represents the local measurements of the signal x made by j th sensor and thus $y = \Psi x$ represents the total measurements of x captured by the whole network. Or, as mentioned previously, DBD matrices can be used to represent the process of measuring a video sequence frame-by-frame, but where each frame is observed using a different measurement matrix. DBD matrices also arise in the study of observability matrices in certain linear dynamical systems [29], and as another example, DBD matrices can be used as a simplified model for the visual pathway (because the information captured by the photoreceptors in the retina is aggregated locally by horizontal and bipolar cells [15]). Due to their structure, DBD matrices can be transformed into sparse measurement matrices after permutation of their rows and columns [1].

We have previously derived concentration inequalities for DBD matrices populated with sub-Gaussian random variables [28, 20]. Rather than ensuring the stable embedding of an entire family of sparse signals, these equalities concern the probability that a bound such as (2) will hold for a single, arbitrary (not necessarily sparse) signal x . We have shown that, unlike the case for unstructured random matrices, the probability of concentration with DBD matrices is actually *signal dependent*, and in particular the concentration probability depends on the allocation of the signal energy among the signal blocks. However, for signals whose energy is nearly uniformly spread across the J blocks (this happens, for example, with signals that are sparse in the Fourier domain [20]), the highly structured DBD matrices can provide concentration performance that is on par with the unstructured matrices often used in CS.

While concentration of measure inequalities are useful for applications concerning compressive signal processing [10], it is not evident how such a concentration result can be extended to give an RIP bound as strong as the one in this paper. Specifically, in Section 3.2 of this paper, we show that if the total number of measurements \bar{M} scales linearly with the sparsity of the signal S and poly-logarithmically with the ambient dimension \tilde{N} , DBD matrices populated with sub-Gaussian random entries will satisfy the U -RIP with high probability. In addition to a dependence on S and \tilde{N} , however, our measurement bounds also reveal a dependence on a property known as the *coherence* of the sparsifying basis U . In this sense, the signal-dependent nature of our concentration of measure inequalities carries over to our RIP analysis for DBD matrices. (The fine details of how this occurs, however, are different.) Our study does confirm that for the class of signals that are sparse in the frequency domain, DBD matrices satisfy the RIP with approximately the same number of rows required in an unstructured Gaussian random matrix (despite having many fewer nonzero entries).

1.3 The RIP for Repeated Block Diagonal (RBD) Matrices

Alternatively, suppose the matrices $\{\Phi_j\}$, $j \in [J]$, in (1) are all equal, and let Ξ denote the resulting block diagonal matrix. Following [20], we say that Ξ has an RBD structure. In the context of the sensor network, video processing, and observability applications discussed before, RBD matrices arise when the same measurement matrix is used for all the signal blocks. In the delay embedding of dynamical systems, as another example, a time series is obtained by repeatedly applying a scalar measurement function to the trajectory of a dynamical system. This time series can then

be embedded in a low-dimensional space (hence the name), and this embedding can be expressed through an RBD measurement matrix provided that the scalar measurement function is linear [32]. Though not obvious at first glance, RBD matrices also have structural similarities with random convolution matrices found in the CS literature [22, 16]. We revisit this connection in order to re-derive the RIP for partially circulant random matrices in Section 3.3.2.

We have previously derived concentration inequalities for RBD matrices populated with Gaussian random variables [28, 23, 20]. Our bounds for these matrices again reveal that the probability of concentration is signal dependent. However, in this case, our bounds depend on both the allocation of the signal energy among the signal blocks as well as the mutual orthogonality of the signal blocks.

In Section 3.3 of this paper, we show that if the total number of measurements \widetilde{M} scales linearly with S and poly-logarithmically with \widetilde{N} , RBD matrices populated with sub-Gaussian random variables will satisfy the U -RIP with high probability. Our measurement bounds also reveal a dependence on a property known as the *block-coherence* of the sparsifying basis U that quantifies the dependence between its row blocks. When the block-coherence of U is small, RBD matrices perform favorably compared to unstructured Gaussian random matrices. Most sparsifying bases are in fact favorable in this regard; we prove that the block-coherence of U is small when U is selected randomly. Once again, for RBD matrices, the signal dependent nature of concentration inequalities and the dependence of the RIP on the sparsifying basis emerge as two sides of the same coin.

1.4 Outline

This paper is organized as follows. Section 2 introduces the notation used throughout the rest of the paper. Section 3 summarizes our main results regarding the RIP for DBD and RBD matrices; these results are later proved in Section 5. Section 4 presents numerical simulations that illustrate the dependence of signal recovery performance on the sparsifying basis U . We conclude the paper with a short discussion in Section 6. We note that, for the reader's convenience, the Toolbox (A) gathers some general tools from linear algebra and probability theory used in our analysis.

2 Notation

We reserve the letters C, C_1, C_2, \dots to represent universal positive constants. We adopt the following (semi-)order: $a \lesssim b$ means that there is an absolute constant C_1 such that $a \leq C_1 b$. If the constant depends on some parameter c , we write $a \lesssim_c b$. Also $a \gtrsim b$ and $a \gtrsim_c b$ are defined similarly.

For an integer S , a signal with no more than S nonzero entries is called S -sparse, and S is known as the sparsity level. In particular, $\|a\|_0$ denotes the number of nonzero entries of a vector a . More generally, a signal that is a linear combination of at most S columns of a basis is said to be S -sparse in that basis. The conjugate transpose of a matrix A will be denoted by A^* . In this paper, $\text{Rank}(A)$ stands for the rank of matrix A . In addition to the regular ℓ_p -norms in the Euclidean spaces, $1 \leq p \leq \infty$, we use $\|A\|_2$ and $\|A\|_F$ to denote the spectral and Frobenius norms of a matrix A , respectively. We use $\|A\|_{\max}$ to denote the largest entry of the matrix A in magnitude. For $1 \leq p \leq \infty$, the Schatten norm of order p of a matrix A is denoted by $\|A\|_{S_p}$ and is defined as

$$\|A\|_{S_p} := \|\sigma_A\|_p,$$

where σ_A is the vector formed by the singular values of A . Observe that $\|A\|_{S_\infty} = \|A\|_2$ and $\|A\|_{S_2} = \|A\|_F$. Throughout this paper, for a matrix A , $\text{vec}(A)$ returns the vector formed by stacking the columns of A . Also, we will use the conventions $[N] := \{1, 2, \dots, N\}$ (for an integer N) and $\#T$ for the cardinality of a set T .

When it appears, the subscript of an expectation operator \mathbb{E} specifies the (group of) random variable(s) with respect to which the expectation is taken. For a random variable Z taking values in \mathbb{C} , we define $\mathbb{E}^p|Z| := (\mathbb{E}|Z|^p)^{1/p}$, $p \geq 1$. A random variable Z is sub-Gaussian if its sub-Gaussian norm, defined below, is finite [27]:

$$\|Z\|_{\psi_2} := \sup_{p \geq 1} \frac{1}{\sqrt{p}} \mathbb{E}^p|Z|. \quad (4)$$

Qualitatively speaking, the tail of (the distribution of) a sub-Gaussian random variable is similar to that of a Gaussian random variable, hence the name. Finally, a Rademacher sequence is a sequence of i.i.d. random variables that take the values ± 1 with equal probability (and are independent of everything else in their every appearance in this paper). In this paper, $\stackrel{\text{i.d.}}{=}$ means that the random variables on both sides of the equality have the same distribution.

A set $\mathcal{C}(\mathcal{S}, \|\cdot\|, r)$ is called a cover for the set \mathcal{S} at resolution r and with respect to the metric $\|\cdot\|$ if for every $x \in \mathcal{S}$, there exists $x' \in \mathcal{C}(\mathcal{S}, \|\cdot\|, r)$ such that $\|x - x'\| \leq r$. The minimum cardinality of all such covers is called the covering number of \mathcal{S} at resolution r and with respect to the norm $\|\cdot\|$, and is denoted here by $\mathcal{N}(\mathcal{S}, \|\cdot\|, r)$.

3 Main Results

3.1 Measures of Coherence

Our results for random block diagonal matrices depend on certain properties of the sparsity basis, i.e., the basis in which the signals have a sparse expansion. These properties are defined and studied in this section; this sets the stage for a detailed statement of our main results in Sections 3.2 and 3.3.

3.1.1 Coherence Definitions

The *coherence* of an orthobasis $U \in \mathbb{C}^{\tilde{N} \times \tilde{N}}$ is defined as follows [7]:

$$\mu(U) := \sqrt{\tilde{N}} \max_{p, q \in [\tilde{N}]} |U(p, q)|, \quad (5)$$

where $U(p, q)$ is the (p, q) th entry of U . If $\{u_{\tilde{n}}\}$ and $\{e_{\tilde{n}}\}$, $\tilde{n} \in [\tilde{N}]$, denote the columns of U and of the canonical basis for $\mathbb{C}^{\tilde{N}}$, respectively, one can easily verify that

$$\mu(U) = \sqrt{\tilde{N}} \max_{p, q \in [\tilde{N}]} |\langle u_p, e_q \rangle|. \quad (6)$$

This allows us to interpret $\mu(U)$ as the similarity between U and the canonical basis.

A few more definitions are in order before we can define the second important property of a basis used in this paper. For $\alpha \in \mathbb{C}^{\tilde{N}}$, set $x(\alpha) = x(\alpha, U) := U\alpha$, and define $x_j(\alpha) = x_j(\alpha, U) \in \mathbb{C}^N$, $j \in [J]$, such that

$$x(\alpha) = [x_1(\alpha)^T, x_2(\alpha)^T, \dots, x_J(\alpha)^T]^T. \quad (7)$$

If we also define $U_j \in \mathbb{C}^{N \times \tilde{N}}$, $j \in [J]$, such that

$$U = [U_1^T, U_2^T, \dots, U_J^T]^T, \quad (8)$$

we observe that $x_j(\alpha) = U_j \alpha$ for every j . Define $X_R(\alpha, U) \in \mathbb{C}^{N \times J}$ as

$$X_R(\alpha) = X_R(\alpha, U) := \begin{bmatrix} x_1(\alpha) & x_2(\alpha) & \cdots & x_J(\alpha) \end{bmatrix} = \begin{bmatrix} U_1 \alpha & \cdots & U_J \alpha \end{bmatrix}.$$

Now the *block-coherence* of U , denoted by $\gamma(U)$, is defined as

$$\gamma(U) := \sqrt{J} \max_{\tilde{n} \in [\tilde{N}]} \|X_R(e_{\tilde{n}}, U)\|_2. \quad (9)$$

In words, $\gamma(U)$ is proportional to the maximal spectral norm when any column of U is reshaped into an $N \times J$ matrix. In analogy with (6), one can also think of (9) as a (non-commutative) coherence measure between U and $I_{\tilde{N}}$.² Qualitatively speaking, $\gamma(U)$ measures the orthogonality and distribution of energy between the row-blocks of U . If for every column of U , the energy is evenly distributed between its row-blocks and they are nearly orthogonal, $\gamma(U)$ will be small and, as we will see later, better suited for our purposes. In the next subsection, we compute the coherence and block-coherence of a few widely-used orthobases.

3.1.2 Computing the Coherence for a Few Orthonormal Bases

It is easily verified that

$$1 \leq \mu(U) \leq \sqrt{\tilde{N}}. \quad (10)$$

The upper bound is achieved, for example, by the canonical basis in $\mathbb{C}^{\tilde{N}}$, i.e., $\mu(I_{\tilde{N}}) = \sqrt{\tilde{N}}$. The lower bound, on the other hand, is achieved by any basis that is maximally incoherent with the canonical basis. For example, $\mu(F_{\tilde{N}}) = 1$, where $F_{\tilde{N}}$ denotes the Fourier basis in $\mathbb{C}^{\tilde{N}}$. The next lemma, proved in B, indicates that *most* orthobases are also highly incoherent with the canonical basis.

Lemma 1. *Let $R \in \mathbb{R}^{\tilde{N} \times \tilde{N}}$ denote a generic orthobasis in $\mathbb{C}^{\tilde{N}}$ chosen randomly from the uniform distribution on the orthogonal group. Then the following holds for fixed $t \gtrsim 1$ and $\tilde{N} \gtrsim t^2 \log \tilde{N}$:*

$$\mathbb{P}\left\{\mu(R) > t\sqrt{\log \tilde{N}}\right\} \lesssim \tilde{N}^{-t}. \quad (11)$$

We now turn to computing the block-coherence of the same orthobases. Since every column of U has unit ℓ_2 -norm, it is easily observed that

$$1 \leq \gamma(U) \leq \sqrt{J}. \quad (12)$$

Consider the canonical basis in $\mathbb{C}^{\tilde{N}}$. For every $\tilde{n} \in [\tilde{N}]$, $X(e_{\tilde{n}}, I_{\tilde{N}})$ has a single non-zero entry, which equals 1, and thus $\|X(e_{\tilde{n}}, I_{\tilde{N}})\|_2 = 1$. Hence, $\gamma(I_{\tilde{N}}) = \sqrt{J}$. Moving on to $F_{\tilde{N}}$, we observe that the entries of the first column of $F_{\tilde{N}}$ equal $\tilde{N}^{-1/2}$. As a result, the entries of $X(e_1, F_{\tilde{N}})$ all equal $\tilde{N}^{-1/2}$. It follows that $\|X(e_1, F_{\tilde{N}})\|_2 = 1$ and therefore $\gamma(F_{\tilde{N}}) = \sqrt{J}$.

²It can be easily verified that, in general, $\max_{\tilde{n}} \|X_R(e_{\tilde{n}}, U)\|_2 \neq \max_{\tilde{n}} \|X_R(u_{\tilde{n}}, I_{\tilde{N}})\|_2$, where $u_{\tilde{n}}$ is the \tilde{n} th column of U .

Because the canonical basis and the Fourier basis—which one might naturally consider to be opposite ends on some spectrum of orthobases—both have large block-coherence, one might wonder whether any orthobasis could have small block-coherence. As we will see, *most* possible orthobases actually do have small block-coherence. For example, consider the generic orthobasis R constructed in Lemma 1. The columns of $X(e_1, R)$ are J random vectors in \mathbb{R}^N . These vectors are weakly dependent because the first column of R has unit ℓ_2 -norm. With high probability, the length of each vector is approximately $1/\sqrt{J}$ (so that the ℓ_2 -norm of the first column of R is one). If $J \leq N$, then with high probability these points are spread out in \mathbb{R}^N so that $\|X(e_1, R)\|_2 \approx 1/\sqrt{J}$. Now, since the columns of R have the same distribution, $\|X(e_{\tilde{n}}, R)\|_2 \approx 1/\sqrt{J}$ for every $\tilde{n} \in [\tilde{N}]$. Therefore, $\gamma(R) \approx 1$, which is much smaller than the block-coherence of the canonical and Fourier bases. The next result, proved in C, formalizes this discussion.

Lemma 2. *Consider the generic orthobasis R constructed in Lemma 1. For fixed $t \leq 1$, the following holds if $J \leq N$ and $N \gtrsim t^{-2} \log \tilde{N}$:*

$$\mathbb{P}\left\{\gamma(R) \gtrsim 1 + \sqrt{\frac{J}{N}} + t\right\} \lesssim \tilde{N}^{-t}. \quad (13)$$

We close this section by noting that Section 3.3.2 provides an example of a deterministic basis with small block-coherence. (This is then used to prove the RIP for partial random circulant matrices.)

3.2 The RIP for DBD Matrices

Let $\Phi \in \mathbb{R}^{M \times N}$ denote a matrix populated with i.i.d. sub-Gaussian random variables having mean zero, standard deviation $1/\sqrt{M}$, and sub-Gaussian norm τ/\sqrt{M} , for some $\tau > 0$. Take $\{\Phi_j\}$ in (1) to be J independent copies of Φ and let $\Psi \in \mathbb{R}^{\tilde{M} \times \tilde{N}}$ denote the resulting block diagonal matrix in (1). Our first main result, proved in Section 5, establishes the RIP for DBD matrices with this construction.

Theorem 1. *Let U denote an orthobasis for $\mathbb{C}^{\tilde{N}}$ and define $\tilde{\mu}(U) := \min(\sqrt{J}, \mu(U))$. If $S \gtrsim 1$ and*

$$\tilde{M} \gtrsim_{\tau} \delta^{-2} \tilde{\mu}^2(U) \cdot S \cdot \log^2 S \log^2 \tilde{N}, \quad (14)$$

then $\delta_S(\Psi, U) \leq \delta < 1$, except with a probability of at most $O(\tilde{N}^{-\log \tilde{N} \log^2 S})$.

A few remarks are in order. The requisite number of measurements is linear in the sparsity level S and (poly-)logarithmic in the ambient dimension \tilde{N} , on par with an unstructured random Gaussian matrix [4]. More importantly, the requisite number of measurements scales with $\tilde{\mu}^2(U)$ which takes a value in the interval $[1, J]$. For the Fourier basis, we calculated that $\mu(F_{\tilde{N}}) = 1$. Therefore, when measuring signals that are sparse in the frequency domain, we observe that a DBD matrix compares favorably to an unstructured Gaussian matrix of the same size. This is in the sense that they both require the same number of measurements to achieve the RIP (up to a poly-logarithmic factor).

On the other hand, when the orthobasis U is highly coherent with the canonical basis, the requisite number of measurements is proportional to SJ (instead of S). While possibly unfavorable, this is indeed necessary (to within a poly-logarithmic factor) to achieve the RIP in some cases. For

example, recall that $\mu(I_{\tilde{N}}) = \sqrt{\tilde{N}}$ and so $\tilde{\mu}(I_{\tilde{N}}) = \sqrt{J}$. To see why the results are optimal in this case, consider the class of S -sparse signals in $I_{\tilde{N}}$ whose nonzero entries are located within the first length- N block of the signal. Achieving a stable embedding of this class of signals requires Φ_1 itself to satisfy the RIP. This matrix, Φ_1 , is an unstructured sub-Gaussian matrix, and ensuring that it satisfies the RIP requires $M \gtrsim_{\tau} \delta^{-2} S \log(N/S)$ [4]. Consequently, ensuring the $I_{\tilde{N}}$ -RIP for Ψ is only possible when $\tilde{M} \gtrsim_{\tau} \delta^{-2} J S \log(N/S)$, as predicted by Theorem 1 (up to a poly-logarithmic term). The required number of measurements in this case can still be parsimonious ($\tilde{M} \ll \tilde{N}$), however, if the sparsity level S of the signal x is much less than N , the length of each signal block x_j .

As a final note, Theorem 1 implies the RIP for a certain class of sparse matrices which are of potential interest in their own right [3, 1].

Corollary 1. *Let U denote an orthobasis for $\mathbb{C}^{\tilde{N}}$ and define $\tilde{\mu}(U) := \min(\sqrt{J}, \mu(U))$. Let Ψ' denote the (sparse) matrix obtained by an arbitrary permutation of the rows and columns of Ψ . If $S \gtrsim 1$ and*

$$\tilde{M} \gtrsim_{\tau} \delta^{-2} \tilde{\mu}^2(U) \cdot S \cdot \log^2 S \log^2 \tilde{N}, \quad (15)$$

then $\delta_S(\Psi', U) \leq \delta < 1$, except with a probability of at most $O(\tilde{N}^{-\log \tilde{N} \log^2 S})$.

Proof. Without loss of generality, consider no permutation in the rows and let $P_c \in \mathbb{R}^{\tilde{N}}$ denote the permutation matrix for the columns of Ψ . Since $P_c U$ has the same coherence as U , the claim follows by applying Theorem 1. \square

3.3 The RIP for RBD Matrices

3.3.1 Main Result for RBD Matrices

Let $\Phi \in \mathbb{R}^{M \times N}$ denote a matrix populated with i.i.d. sub-Gaussian random variables having mean zero, standard deviation $1/\sqrt{M}$, and sub-Gaussian norm τ/\sqrt{M} , for some $\tau > 0$. Take $\Phi_j = \Phi$ for every $j \in [J]$ in (1) and let $\Xi \in \mathbb{R}^{\tilde{M} \times \tilde{N}}$ denote the resulting block diagonal matrix in (1). Our second main result, also proved in Section 5, establishes the RIP for RBD matrices with this construction.

Theorem 2. *Let U denote an orthobasis for $\mathbb{C}^{\tilde{N}}$. If $S \gtrsim 1$ and*

$$\tilde{M} \gtrsim_{\tau} \delta^{-2} \gamma^2(U) \cdot S \cdot \log^2 S \log^2 \tilde{N},$$

then $\delta_S(\Xi, U) \leq \delta < 1$, except with a probability of at most $O(\tilde{N}^{-\log \tilde{N} \log^2 S})$.

A few remarks are in order. At one end of the spectrum, the block-coherence of an orthobasis could equal \sqrt{J} and consequently the required number of measurements above would scale with JS . This happens for signals that are sparse, for example, in the time (canonical basis) or frequency domains. Our result is indeed optimal for both of these bases (up to a poly-logarithmic factor). The same argument for the canonical basis carries over from the DBD matrices. For the Fourier basis, we note that it is possible to construct certain classes of periodic signals in $\mathbb{C}^{\tilde{N}}$ that would require the lower bound on \tilde{M} to scale with JS . Consider, for example, the class of signals consisting of all S -sparse combinations of columns $1, J+1, \dots, (J-1)N+1$ from $F_{\tilde{N}}$. If x belongs to this class,

then, by construction, x_1, x_2, \dots, x_J (as defined in (7)) are all equal because x is periodic with a period of N . As a result, different blocks of Ξ take the same measurements from x . Therefore, as was the case with the DBD matrices, obtaining a stable embedding of this class of signals requires $M \gtrsim_{\tau} \delta^{-2} S \log(N/S)$, and equivalently, $\bar{M} \gtrsim_{\tau} \delta^{-2} JS \log(N/S)$.

At the other end of the spectrum, for a generic orthobasis R we computed that $\gamma(R) \lesssim 1$ with high probability. For signals that are sparse in this basis (and therefore for many possible orthobases in general), an RBD matrix performs nearly as well as an unstructured Gaussian random matrix. The RBD structure allows us to prove the RIP for certain classes of structured random matrices as a special case of Theorem 2. In particular, in the next subsection we re-derive an RIP bound for partial random circulant matrices that originally appeared in [16]. As a byproduct, we also construct a deterministic sparsity basis that achieves a performance similar to the generic orthobasis we have considered above.

3.3.2 The RIP for Partial Random Circulant Matrices

This section demonstrates that the RBD model, together with Theorem 2, can be used to derive the RIP for partial random circulant matrices.³ More specifically, we focus on proving the RIP for $\Gamma \in \mathbb{R}^{J \times P}$, with $J \leq P$, defined as

$$\Gamma = \frac{1}{\sqrt{J}} \begin{bmatrix} \epsilon_1 & \epsilon_2 & \cdots & \epsilon_P \\ \epsilon_P & \epsilon_1 & \cdots & \epsilon_{P-1} \\ \vdots & \vdots & & \vdots \\ \epsilon_{P-J+2} & \epsilon_{P-J+3} & \cdots & \epsilon_{P-J+1} \end{bmatrix},$$

where $\{\epsilon_p\}$, $p \in [P]$, is a sequence of i.i.d. zero-mean, unit-variance random variables with sub-Gaussian norm τ . We let ε denote the vector formed by $\{\epsilon_p\}$. In order to use Theorem 2 in this setting, we make the following argument: for any signal $x \in \mathbb{C}^P$, we can write Γx as the multiplication of an RBD matrix $\Xi \in \mathbb{R}^{J \times PJ}$ and an extended vector $\hat{x} \in \mathbb{C}^{PJ}$:

$$\Gamma x = \overbrace{\begin{bmatrix} \varepsilon^* & & & \\ & \varepsilon^* & & \\ & & \ddots & \\ & & & \varepsilon^* \end{bmatrix}}^{\Xi} \cdot \frac{1}{\sqrt{J}} \overbrace{\begin{bmatrix} S^0 x \\ S^1 x \\ \vdots \\ S^{J-1} x \end{bmatrix}}^{\hat{x}} =: \Xi \cdot \frac{1}{\sqrt{J}} \hat{x}, \quad (16)$$

where S is the cyclic shift-up operator on column vectors in \mathbb{C}^P . The next lemma (proved in D) states that if x is sparse, then \hat{x} has a sparse representation in a favorable orthobasis T (constructed in the proof).

Lemma 3. *There exists an orthobasis T with $\gamma(T) = 1$, such that every \hat{x} has an S -sparse representation in T if the corresponding x is S -sparse. That is, for every such \hat{x} , there exists x_e with $\|x_e\|_0 \leq S$ such that $\hat{x}/\sqrt{J} = Tx_e$.*

³The arguments in this section extend without much effort to the more general case of establishing the RIP for partial random Toeplitz matrices.

Therefore, the T -RIP for Ξ implies the RIP for Γ . To be more specific, after setting $M = 1$, Theorem 2 implies that $\delta_S(\Gamma) \leq \delta_S(\Xi, T) \leq \delta$ except with a probability of at most $O(P^{-\log P \log^2 S})$, provided that

$$J \gtrsim_{\tau} \delta^{-2} \cdot S \log^2 S \log^2 P,$$

which is equivalent to Theorem 1.1 in [16].

4 Numerical Simulations

This section contains a series of simulations that are intended to reinforce our findings for the reader. An ideal scenario would involve generating several random block diagonal matrices while varying the sparsity level S and the number of measurements \widetilde{M} and measuring the fraction of realizations in which the RIC falls below a fixed threshold. Checking the RIP for a matrix is, however, known to be an NP hard problem [26]; this encourages us to examine a proxy for the RIP. As discussed in Section 1.1, a major application of the RIP is to ensure robust recovery of sparse signals via algorithms such as Basis Pursuit (BP).⁴ Similar to [13], then, we measure the success of sparse signal recovery as an alternative to verifying the RIP. This is detailed in the following.

With $N = 100$ and $J = 10$ (and consequently, $\widetilde{N} = 1000$), we generate an $\widetilde{N} \times \widetilde{N}$ DBD matrix whose entries are i.i.d. Gaussian random variables having zero mean and unit standard deviation. Note that each diagonal block is of size $N \times N$. For each pair $(S, M) \in [\widetilde{M}] \times [N]$, the following procedure is executed. An $\widetilde{M} \times \widetilde{N}$ DBD matrix Ψ is formed by keeping (and appropriately normalizing) the first M rows of each diagonal block of the large $\widetilde{N} \times \widetilde{N}$ matrix. Then 20 random S -sparse signals in the canonical basis of $\mathbb{C}^{\widetilde{N}}$ are generated. These sparse signals are measured using Ψ and reconstructed back (from incomplete measurements) via BP.⁵ We deem the recovery successful if the relative ℓ_2 reconstruction error is less than the fixed threshold 10^{-2} . The fraction of successful recovery is recorded and this procedure is repeated for other pairs (S, M) . The resulting phase transition graph is depicted in Figure 1a, where the color of each pixel ranges from black for perfect recovery in every realization to white for failed recovery every time.

We repeat the above simulation for signals that are sparse in the Fourier and generic bases (see Lemma 1). A new generic basis is generated in each iteration. All of the above simulations are then repeated with an RBD matrix. The results of the simulations are displayed in Figures 1 and 2. In the simulations with DBD matrices, recovery of canonical and frequency sparse signals are, respectively, the least and most successful of the three instances. Signals that are sparse in the generic bases can be recovered nearly as well as frequency sparse signals.

In the simulations with RBD matrices, signals that are sparse in the generic bases are recovered best, while the results for the canonical and Fourier bases are less satisfactory. These observations are in agreement with our findings in Section 3. We do point out, however, that for the case of RBD measurement matrices we are able to recover frequency sparse signals somewhat better than signals that are sparse in the canonical basis. We note that such difference is not reflected in our RIP bounds. While this performance difference could be due to a lack of explicit numerical constants in our results, it is more likely an artifact of our simulations: We are not directly confirming the RIP but rather testing the recovery of randomly generated test signals, and those few signals which

⁴In a nutshell, BP casts the signal recovery problem as a linear program, which can be efficiently solved. We refer the reader to [5] for more detail.

⁵In order to implement BP, we used YALL1, a package for solving ℓ_1 problems [34, 30].

make it difficult to satisfy the RIP may be more pathological in the Fourier basis than in the canonical basis.

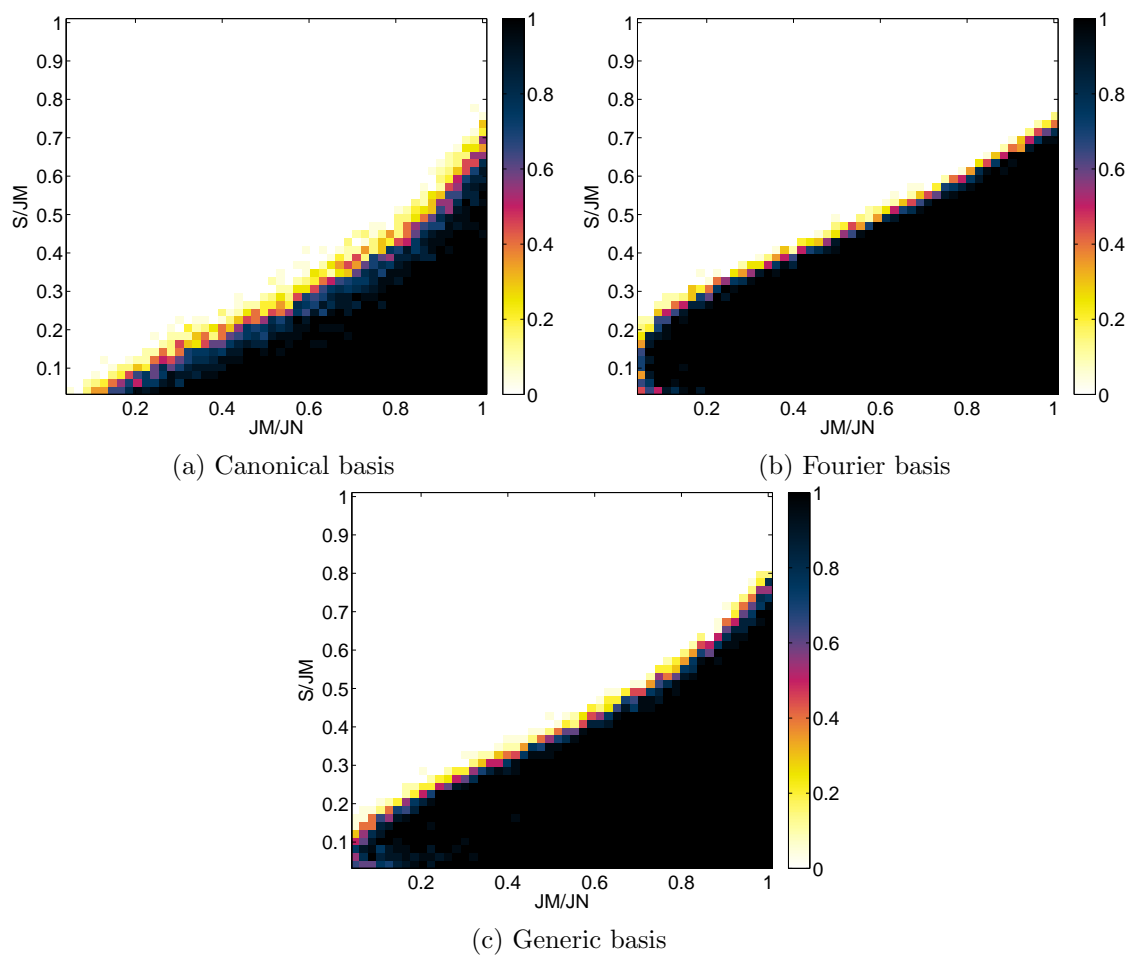


Figure 1: Simulation results for DBD matrices. Refer to Section 4 for details.

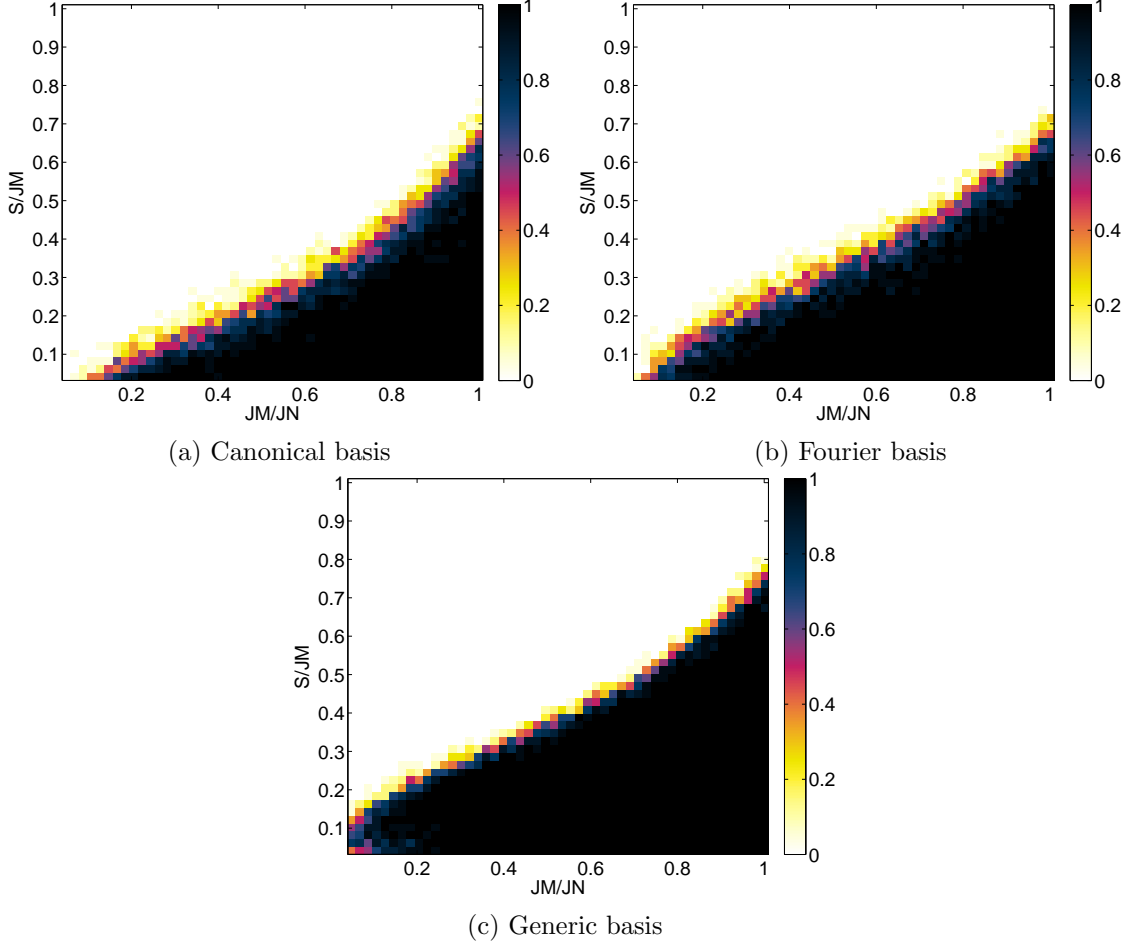


Figure 2: Simulation results for RBD matrices. Refer to Section 4 for details.

5 Proofs of Theorems 1 and 2

5.1 Preliminaries

First, define the set of all S -sparse signals with unit norm as

$$\Omega_S := \left\{ \alpha \in \mathbb{C}^{\tilde{N}} : \|\alpha\|_0 \leq S, \|\alpha\|_2 = 1 \right\}. \quad (17)$$

With this definition, we observe that the RIC for DBD matrices can be written as

$$\delta_S = \sup_{\alpha \in \Omega_S} \left| \|\Psi \cdot x(\alpha)\|_2^2 - 1 \right|,$$

where we leveraged the fact that $\|\alpha\|_2 = 1$ implies

$$\mathbb{E}\{\|\Psi \cdot x(\alpha)\|_2^2\} = \alpha^* U^* \mathbb{E}\{\Psi^* \Psi\} U \alpha = \alpha^* \alpha = 1.$$

The RIC for RBD matrices can be similarly written as

$$\delta_S = \sup_{\alpha \in \Omega_S} \left| \|\Xi \cdot x(\alpha)\|_2^2 - 1 \right|.$$

Given $\delta < 1$ and under the conditions in Theorems 1 and 2, our objective is to show that $\delta_S \leq \delta$ for both DBD and RBD matrices. To achieve this goal, we require the following powerful result due to Krahmer et al.:

Theorem 3. [16, Theorem 3.1] *Let $\mathcal{A} \subset \mathbb{C}^{\widetilde{M} \times \widetilde{N}}$ be a set of matrices, and let ε be a random vector whose entries are i.i.d., zero-mean, unit-variance random variables with sub-Gaussian norm τ . Set*

$$\begin{aligned} d_F(\mathcal{A}) &:= \sup_{A \in \mathcal{A}} \|A\|_F, \\ d_2(\mathcal{A}) &:= \sup_{A \in \mathcal{A}} \|A\|_2, \end{aligned}$$

and

$$\begin{aligned} E_1 &:= \gamma_2(\mathcal{A}, \|\cdot\|_2) (\gamma_2(\mathcal{A}, \|\cdot\|_2) + d_F(\mathcal{A})) + d_F(\mathcal{A}) d_2(\mathcal{A}), \\ E_2 &:= d_2(\mathcal{A}) (\gamma_2(\mathcal{A}, \|\cdot\|_2) + d_F(\mathcal{A})), \\ E_3 &:= d_2^2(\mathcal{A}). \end{aligned}$$

Then, for $t > 0$, it holds that

$$\log P \left\{ \sup_{A \in \mathcal{A}} \left| \|A\varepsilon\|_2^2 - \mathbb{E} \|A\varepsilon\|_2^2 \right| \gtrsim_\tau E_1 + t \right\} \lesssim_\tau - \min \left(\frac{t^2}{E_2^2}, \frac{t}{E_3} \right).$$

Without going into the details, we note that the γ_2 -functional of \mathcal{A} , $\gamma_2(\mathcal{A}, \|\cdot\|_2)$, is a geometrical property of \mathcal{A} , i.e., the index set of the random process, and is widely used in the context of probability in Banach spaces [25, 18]. In particular, the following lemma gives an estimate of this quantity.

Lemma 4. [25, 17] *With \mathcal{A} as defined above, it holds that*

$$\gamma_2(\mathcal{A}, \|\cdot\|_2) \lesssim \int_0^\infty \log^{\frac{1}{2}}(\mathcal{N}(\mathcal{A}, \|\cdot\|_2, \nu)) \, d\nu. \quad (18)$$

Clearly, we need to express the problem of bounding the RIC of DBD and RBD matrices in a form that is amenable to the setting of Theorem 3. First, for DBD matrices, let us define $X_{D,j} \in \mathbb{C}^{M \times MN}$, $j \in [J]$, as

$$X_{D,j}(\alpha) = X_{D,j}(\alpha, U) := \begin{bmatrix} x_j^*(\alpha) & & & \\ & x_j^*(\alpha) & & \\ & & \ddots & \\ & & & x_j^*(\alpha) \end{bmatrix}. \quad (19)$$

It can then be easily verified that

$$\begin{aligned} \|\Psi \cdot x(\alpha)\|_2^2 &= \sum_{j \in [J]} \|\Phi_j \cdot x_j(\alpha)\|_2^2 \\ &= \sum_{j \in [J]} \|X_{D,j}(\alpha) \cdot \text{vec}(\Phi_j^*)\|_2^2 \\ &\stackrel{\text{i.d.}}{=} \sum_{j \in [J]} \left\| \frac{1}{\sqrt{M}} X_{D,j}(\alpha) \cdot \varepsilon_j \right\|_2^2 \\ &\stackrel{\text{i.d.}}{=} \|A_D(\alpha) \cdot \varepsilon\|_2^2, \end{aligned}$$

where the linear map $A_D : \Omega_S \rightarrow \mathbb{C}^{\widetilde{M} \times JMN}$ is defined as

$$A_D(\alpha) = A_D(\alpha, U) := \frac{1}{\sqrt{M}} \begin{bmatrix} X_{D,1}(\alpha) & & & \\ & X_{D,2}(\alpha) & & \\ & & \ddots & \\ & & & X_{D,J}(\alpha) \end{bmatrix},$$

and entries of $\varepsilon_j \in \mathbb{R}^{MN}$, $j \in [J]$, and $\varepsilon \in \mathbb{R}^{JMN}$ are i.i.d. zero-mean, unit-variance random variables with sub-Gaussian norm τ . The index set of the random process is $\mathcal{A}_D := \{A_D(\alpha) : \alpha \in \Omega_S\}$. We have therefore completely expressed the DBD problem in the setting of Theorem 3.

Next, for RBD matrices, we observe that

$$\begin{aligned} \|\Xi \cdot x(\alpha)\|_2^2 &= \sum_{j \in [J]} \|\Phi \cdot x_j(\alpha)\|_2^2 \\ &= \|\Phi \cdot X_R(\alpha)\|_F^2 \\ &= \|X_R^*(\alpha) \cdot \Phi^*\|_F^2 \\ &\stackrel{\text{i.d.}}{=} \sum_{m \in [M]} \left\| \frac{1}{\sqrt{M}} X_R^*(\alpha) \cdot \varepsilon'_m \right\|_2^2 \\ &\stackrel{\text{i.d.}}{=} \|A_R(\alpha) \cdot \varepsilon'\|_2^2, \end{aligned} \tag{20}$$

where in the last line we defined the linear map $A_R : \Omega_S \rightarrow \mathbb{C}^{\widetilde{M} \times MN}$ as

$$A_R(\alpha) = A_R(\alpha, U) := \frac{1}{\sqrt{M}} \begin{bmatrix} X_R^*(\alpha) & & & \\ & X_R^*(\alpha) & & \\ & & \ddots & \\ & & & X_R^*(\alpha) \end{bmatrix}, \tag{21}$$

and the entries of $\varepsilon'_m \in \mathbb{R}^N$, $m \in [M]$, and $\varepsilon' \in \mathbb{R}^{MN}$ are i.i.d. zero-mean, unit-variance random variables with sub-Gaussian norm of τ . So, we have also managed to express the RBD problem in the setting of Theorem 3. The next two subsections are concerned with estimating the quantities involved in Theorem 3 for both the DBD and RBD problems.

5.2 Calculating $d_2(\mathcal{A}_D)$, $d_F(\mathcal{A}_D)$, and $\gamma_2(\mathcal{A}_D, \|\cdot\|_2)$

We begin with defining the following norm on $\mathbb{C}^{\widetilde{N}}$, which will find extensive use in the analysis of DBD problem:

$$\|\alpha\|_{A_D} := \|A_D(\alpha)\|_2 \tag{22}$$

for $\alpha \in \mathbb{C}^{\widetilde{N}}$. We record a useful property of this norm below.

Lemma 5. *For every $\alpha \in \mathbb{C}^{\widetilde{N}}$, it holds that*

$$\|\alpha\|_{A_D} \leq \frac{\widetilde{\mu}}{\sqrt{M}} \cdot \|\alpha\|_1. \tag{23}$$

Proof. Let $u_{j,n}$, $j \in [J]$ and $n \in [N]$, denote the $((j-1)N+n)$ th row of U . We then have that

$$\begin{aligned}
\|\alpha\|_{A_D} &= \|A_D(\alpha)\|_2 \\
&= \|A_D(\alpha)A_D^*(\alpha)\|_2^{\frac{1}{2}} \\
&= \frac{1}{\sqrt{M}} \max_{j \in [J]} \|x_j\|_2 \\
&= \frac{1}{\sqrt{M}} \max_{j \in [J]} \|U_j \alpha\|_2 \\
&\leq \sqrt{\frac{N}{M}} \max_{j \in [J], n \in [N]} |\langle u_{j,n}, \alpha \rangle| \\
&\leq \sqrt{\frac{N}{M}} \max_{j \in [J], n \in [N]} \|u_{j,n}\|_\infty \|\alpha\|_1 \\
&= \frac{\mu}{\sqrt{M}} \cdot \|\alpha\|_1,
\end{aligned}$$

where the second to last line uses the Hölder inequality and the last line follows from the definition of μ . On the other hand, one may also write that

$$\|\alpha\|_{A_D} = \frac{1}{\sqrt{M}} \max_{j \in [J]} \|x_j\|_2 \leq \frac{1}{\sqrt{M}} \|x\|_2 = \frac{1}{\sqrt{M}} \|U\alpha\|_2 = \frac{1}{\sqrt{M}} \|\alpha\|_2 \leq \frac{1}{\sqrt{M}} \|\alpha\|_1,$$

where we used the fact that U is an orthobasis. Overall, we arrive at

$$\|\alpha\|_{A_D} \leq \frac{1}{\sqrt{M}} \min(\mu, \sqrt{J}) \|\alpha\|_1 = \frac{\tilde{\mu}}{\sqrt{M}} \cdot \|\alpha\|_1, \quad (24)$$

as claimed. The equality above follows from the definition of $\tilde{\mu}$. \square

We continue with computing the quantities involved in Theorem 3 in the case of the DBD problem. First, we have that

$$d_F(\mathcal{A}_D) = \sup_{A_D(\alpha) \in \mathcal{A}_D} \|A_D(\alpha)\|_F = \sup_{\alpha \in \Omega_S} \|x(\alpha)\|_2 = \sup_{\alpha \in \Omega_S} \|U\alpha\|_2 = \sup_{\alpha \in \Omega_S} \|\alpha\|_2 = 1. \quad (25)$$

The second to last equality holds because U is an orthonormal basis. Second, we have that

$$d_2(\mathcal{A}_D) = \sup_{A_D(\alpha) \in \mathcal{A}_D} \|A_D(\alpha)\|_2 = \sup_{\alpha \in \Omega_S} \|\alpha\|_{A_D} \leq \frac{\tilde{\mu}}{\sqrt{M}} \sup_{\alpha \in \Omega_S} \|\alpha\|_1 \leq \tilde{\mu} \sqrt{\frac{S}{M}}. \quad (26)$$

The first inequality above holds on account of Lemma 5. The second inequality above follows because $\|\alpha\|_2 = 1$ and $\|\alpha\|_0 \leq S$ when $\alpha \in \Omega_S$. It is only left to bound $\gamma_2(\mathcal{A}_D, \|\cdot\|_2)$. According to Lemma 4, we have that

$$\begin{aligned}
\gamma_2(\mathcal{A}_D, \|\cdot\|_2) &\leq \int_0^\infty \log^{\frac{1}{2}}(\mathcal{N}(\mathcal{A}_D, \|\cdot\|_2, \nu)) \, d\nu \\
&= \int_0^\infty \log^{\frac{1}{2}}(\mathcal{N}(\Omega_S, \|\cdot\|_{A_D}, \nu)) \, d\nu,
\end{aligned}$$

where the isometry between \mathcal{A}_D (with metric $\|\cdot\|_2$) and Ω_S (with metric $\|\cdot\|_{A_D}$) implies the second line. This isometry, in turn, follows from (22) and the linearity of $A_D(\cdot)$. Consequently,

$$\begin{aligned}\gamma_2(\mathcal{A}_D, \|\cdot\|_2) &\leq \int_0^\infty \log^{\frac{1}{2}} \left(\mathcal{N} \left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_{A_D}, \frac{\nu}{\sqrt{S}} \right) \right) d\nu \\ &\leq \sqrt{S} \int_0^\infty \log^{\frac{1}{2}} \left(\mathcal{N} \left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_{A_D}, \nu \right) \right) d\nu,\end{aligned}\tag{27}$$

where the first line uses the second inequality in (47) and the last line follows from a change of variables in the integral. An estimate of the covering number involved in (27) can be found through the next result, which is proved in E.

Lemma 6. *Consider a norm $\|\cdot\|_A$ on $\mathbb{C}^{\tilde{N}}$ that, for every $\alpha \in \mathbb{C}^{\tilde{N}}$, satisfies*

$$\|\alpha\|_A = \|A(\alpha)\|_2 \leq \frac{\kappa}{\sqrt{\widetilde{M}}} \cdot \|\alpha\|_1,$$

for some linear map $A(\cdot) : \mathbb{C}^{\tilde{N}} \rightarrow \mathbb{C}^{N'}$ with rank of at most \widetilde{M} and some $\kappa > 0$ and integer N' . Then, for $0 < \nu < \kappa/\sqrt{\widetilde{M}}$ and $\widetilde{M} \gtrsim 1$, we have that

$$\log \left(\mathcal{N} \left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_A, \nu \right) \right) \lesssim \min \left(S \log \tilde{N} + S \log \left(1 + \frac{2\kappa}{\nu\sqrt{\widetilde{M}}} \right), \frac{\kappa^2}{\nu^2\widetilde{M}} \cdot \log^2 \tilde{N} \right).\tag{28}$$

When $\nu \geq \kappa/\sqrt{\widetilde{M}}$, we have $\mathcal{N} \left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_A, \nu \right) = 1$.

Qualitatively speaking, of the two bounds on the right hand of (28), the first is tighter when ν is small while the second is more effective for larger values of ν . Of course, $\|\cdot\|_{A_D}$ satisfies the hypothesis of Lemma 6 with $\kappa = \tilde{\mu}$ and the map $A_D(\cdot)$. Consequently, for $0 < \nu_0 \leq \tilde{\mu}/\sqrt{\widetilde{M}}$ to be set later, we have that

$$\begin{aligned}&\int_0^\infty \log^{\frac{1}{2}} \left(\mathcal{N} \left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_{A_D}, \nu \right) \right) d\nu \\ &= \int_0^{\nu_0} \log^{\frac{1}{2}} \left(\mathcal{N} \left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_{A_D}, \nu \right) \right) d\nu + \int_{\nu_0}^{\frac{\tilde{\mu}}{\sqrt{\widetilde{M}}}} \log^{\frac{1}{2}} \left(\mathcal{N} \left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_{A_D}, \nu \right) \right) d\nu \\ &\lesssim \int_0^{\nu_0} \left(\sqrt{S \log \tilde{N}} + \sqrt{S \log \left(1 + \frac{2\tilde{\mu}}{\nu\sqrt{\widetilde{M}}} \right)} \right) d\nu + \log \tilde{N} \int_{\nu_0}^{\frac{\tilde{\mu}}{\sqrt{\widetilde{M}}}} \frac{\tilde{\mu}}{\nu\sqrt{\widetilde{M}}} d\nu \\ &\lesssim \nu_0 \sqrt{S \log \tilde{N}} + \nu_0 \sqrt{S \log \left(1 + \frac{2\tilde{\mu}}{\nu_0\sqrt{\widetilde{M}}} \right)} + \frac{\tilde{\mu}}{\sqrt{\widetilde{M}}} \log \tilde{N} \log \left(\frac{\tilde{\mu}}{\nu_0\sqrt{\widetilde{M}}} \right).\end{aligned}\tag{29}$$

The second line above follows from the second statement in Lemma 6. In the third line, different upper bounds from (28) are used to bound each summand. We benefited from (49) in the Toolbox to compute the logarithmic integral in the third line. With the choice of $\nu_0 = \tilde{\mu}/\sqrt{S\widetilde{M}}$, we obtain

that

$$\begin{aligned}
& \int_0^\infty \log^{\frac{1}{2}} \left(\mathcal{N} \left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_{A_D}, \nu \right) \right) d\nu \\
& \lesssim \frac{\tilde{\mu}}{\sqrt{\widetilde{M}}} \sqrt{\log \tilde{N}} + \frac{\tilde{\mu}}{\sqrt{\widetilde{M}}} \sqrt{\log(1 + 2\sqrt{S})} + \frac{\tilde{\mu}}{\sqrt{\widetilde{M}}} \log S \log \tilde{N} \\
& \lesssim \frac{\tilde{\mu}}{\sqrt{\widetilde{M}}} \log S \log \tilde{N}
\end{aligned} \tag{30}$$

for $S \gtrsim 1$. Now plugging back (30) into (27), we arrive at

$$\gamma_2(\mathcal{A}_D, \|\cdot\|_2) \lesssim \tilde{\mu} \sqrt{\frac{S}{\widetilde{M}}} \log S \log \tilde{N}. \tag{31}$$

Before completing the analysis of the DBD problem, let us calculate the same quantities for the RBD case.

5.3 Calculating $d_2(\mathcal{A}_R)$, $d_F(\mathcal{A}_R)$, and $\gamma_2(\mathcal{A}_R, \|\cdot\|_2)$

Again, we first introduce the following norm on $\mathbb{C}^{\tilde{N}}$, which will be useful in the analysis of the RBD problem:

$$\|\alpha\|_{A_R} := \|A_R(\alpha)\|_2 \tag{32}$$

for $\alpha \in \mathbb{C}^{\tilde{N}}$. This norm has the following property.

Lemma 7. *For every $\alpha \in \mathbb{C}^{\tilde{N}}$, it holds that*

$$\|\alpha\|_{A_R} \leq \frac{\gamma}{\sqrt{\widetilde{M}}} \cdot \|\alpha\|_1. \tag{33}$$

Proof. Note that

$$\begin{aligned}
\|\alpha\|_{A_R} &= \|A_R(\alpha)\|_2 \\
&= \frac{1}{\sqrt{\widetilde{M}}} \|X_R(\alpha)\|_2 \\
&= \frac{1}{\sqrt{\widetilde{M}}} \left\| \sum_{\tilde{n} \in [\tilde{N}]} \alpha(\tilde{n}) X_R(e_{\tilde{n}}) \right\|_2 \\
&\leq \frac{1}{\sqrt{\widetilde{M}}} \sum_{\tilde{n} \in [\tilde{N}]} |\alpha(\tilde{n})| \cdot \|X_R(e_{\tilde{n}})\|_2 \\
&\leq \frac{1}{\sqrt{\widetilde{M}}} \max_{\tilde{n} \in [\tilde{N}]} \|X_R(e_{\tilde{n}})\|_2 \cdot \|\alpha\|_1 \\
&= \frac{\gamma}{\sqrt{\widetilde{M}}} \cdot \|\alpha\|_1.
\end{aligned}$$

The third line above uses the linearity of $X_R(\cdot)$. The fourth line follows from the triangle inequality and the fifth line is implied by the Hölder inequality. We made use of the definition of γ to get the last line. \square

We now continue with computing the quantities involved in Theorem 3 in the case of the RBD problem. First, we have that

$$d_F(\mathcal{A}_R) = \sup_{A_R(\alpha) \in \mathcal{A}_R} \|A_R(\alpha)\|_F = \sup_{\alpha \in \Omega_S} \|X_R(\alpha)\|_F = \sup_{\alpha \in \Omega_S} \|x(\alpha)\|_2 = \sup_{\alpha \in \Omega_S} \|\alpha\|_2 = 1. \quad (34)$$

Second, we have that

$$d_2(\mathcal{A}_R) = \sup_{A_R(\alpha) \in \mathcal{A}_R} \|A_R(\alpha)\|_2 = \sup_{\alpha \in \Omega_S} \|\alpha\|_{A_R} \leq \frac{\gamma}{\sqrt{\widetilde{M}}} \sup_{\alpha \in \Omega_S} \|\alpha\|_1 \leq \gamma \sqrt{\frac{S}{\widetilde{M}}}, \quad (35)$$

where the first inequality follows from Lemma 7. The last quantity to estimate is $\gamma_2(\mathcal{A}_R, \|\cdot\|_2)$. As was the case in the previous subsection, we may write that

$$\gamma_2(\mathcal{A}_R, \|\cdot\|_2) \leq \sqrt{S} \int_0^\infty \log^{\frac{1}{2}} \left(\mathcal{N} \left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_{A_R}, \nu \right) \right) d\nu. \quad (36)$$

Of course, $\|\cdot\|_{A_R}$ satisfies the hypothesis of Lemma 6 with $\kappa = \gamma$ and the map $A_R(\cdot)$. Therefore, following the same steps as in the previous subsection, we arrive at

$$\gamma_2(\mathcal{A}_R, \|\cdot\|_2) \lesssim \gamma \sqrt{\frac{S}{\widetilde{M}}} \log S \log \widetilde{N}. \quad (37)$$

5.4 Denouement

We notice that the quantities $d_F(\mathcal{A}_D)$, $d_2(\mathcal{A}_D)$, and $\gamma_2(\mathcal{A}_D, \|\cdot\|_2)$ have the same bounds as their counterparts $d_F(\mathcal{A}_R)$, $d_2(\mathcal{A}_R)$, and $\gamma_2(\mathcal{A}_R, \|\cdot\|_2)$ except for the type of the coherence factor involved. Therefore, it suffices to focus on one—the same result holds for the other with its corresponding coherence factor. Given $\delta < 1$, assume that

$$\widetilde{M} \gtrsim_\tau \delta^{-2} \widetilde{\mu}^2 \cdot S \log^2 S \log^2 \widetilde{N}. \quad (38)$$

Equipped with the estimates in Section 5.2, i.e., (25), (26), and (31), we now compute E_1 in Theorem 3:

$$\begin{aligned} E_1 &:= \gamma_2(\mathcal{A}_D, \|\cdot\|_2) (\gamma_2(\mathcal{A}_D, \|\cdot\|_2) + d_F(\mathcal{A}_D)) + d_F(\mathcal{A}_D) d_2(\mathcal{A}_D) \\ &\lesssim \widetilde{\mu} \sqrt{\frac{S}{\widetilde{M}}} \log S \log \widetilde{N} \left(\widetilde{\mu} \sqrt{\frac{S}{\widetilde{M}}} \log S \log \widetilde{N} + 1 \right) + \widetilde{\mu} \sqrt{\frac{S}{\widetilde{M}}} \\ &\lesssim_\tau \delta (\delta + 1) + \frac{\delta}{\log S \log \widetilde{N}} \\ &\leq 2\delta + \frac{\delta}{\log S \log \widetilde{N}} \\ &\lesssim \delta, \end{aligned}$$

where we assumed that $S \gtrsim 1$ and used the hypothesis that $\delta < 1$ in the last line. Doing the same to E_2 , we obtain that

$$\begin{aligned}
E_2 &:= d_2(\mathcal{A}_D) (\gamma_2(\mathcal{A}_D, \|\cdot\|_2) + d_F(\mathcal{A}_D)) \\
&\lesssim \tilde{\mu} \sqrt{\frac{S}{\widetilde{M}}} \left(\tilde{\mu} \sqrt{\frac{S}{\widetilde{M}}} \log S \log \tilde{N} + 1 \right) \\
&\lesssim_\tau \frac{\delta}{\log S \log \tilde{N}} (\delta + 1) \\
&\lesssim \frac{\delta}{\log S \log \tilde{N}}.
\end{aligned}$$

Similarly for E_3 , we write that

$$E_3 := d_2^2(\mathcal{A}_D) \leq \frac{\tilde{\mu}^2 S}{\widetilde{M}} \lesssim_\tau \frac{\delta^2}{\log^2 S \log^2 \tilde{N}}.$$

Plugging the above estimates of E_1 , E_2 , and E_3 into the tail bound in Theorem 3, we obtain that

$$\log P \left\{ \sup_{\alpha \in \Omega_S} \left| \|\Psi \cdot x(\alpha)\|_2^2 - 1 \right| \gtrsim_\tau \delta + t \right\} \lesssim_\tau - \min \left(\delta^{-2} t^2 \log^2 S \log^2 \tilde{N}, \delta^{-1} t \log^2 S \log^2 \tilde{N} \right).$$

Substituting $t = \delta$, we arrive at

$$\log P \left\{ \sup_{\alpha \in \Omega_S} \left| \|\Psi \cdot x(\alpha)\|_2^2 - 1 \right| \gtrsim_\tau \delta \right\} \lesssim_\tau - \log^2 S \log^2 \tilde{N},$$

assuming that $S \gtrsim 1$. After absorbing the factor depending on τ into (a redefined) δ , we finally arrive at

$$\log P \left\{ \sup_{\alpha \in \Omega_S} \left| \|\Psi \cdot x(\alpha)\|_2^2 - 1 \right| > \delta \right\} \lesssim_\tau - \log^2 S \log^2 \tilde{N},$$

which completes the proof of Theorem 1. Replacing $\tilde{\mu}$ with γ and repeating this argument concludes the proof of Theorem 2.

6 Conclusion

In this paper, we studied two important classes of structured random matrices, namely DBD and RBD matrices. Our main results state that matrices with block diagonal constructions can indeed satisfy the RIP but that the requisite number of measurements depends on certain properties of the sparsifying basis. These properties were detailed and carefully interpreted in the paper. In the best case, DBD and RBD matrices perform nearly as well as the dense i.i.d. random matrices generally used in CS despite having many fewer nonzero entries. Moreover, we have shown that random block diagonal matrices are intimately related to the random convolution and random Toeplitz matrices considered in the literature.

Our findings lead us to conclude that structured random matrices can be useful in sensing architectures as long as the statistics of the data are well-matched with the structure of the measurement matrix. While this intuition is similar to other results on structured measurement matrices, our results on block diagonal matrix constructions are novel in extending this intuition to matrices with (potentially) many entries that are zero. The approach required to reach our results also leads us to conclude that future progress in the field of probability in Banach spaces is likely to play a significant role in establishing optimal performance guarantees for other structured measurement systems. This may be especially true given the improved performance we were able to achieve even over other bounding techniques that require sophisticated mathematical machinery. Finally, while we remain uncertain about the necessity of the poly-logarithmic factors in the final measurement rates (which may be a proof artifact), the simulation results display the dependence on the sparsity basis (through the coherence) present in our main results. Despite the fact that the simulation results address average case behavior (as opposed to worst-case behavior captured by the RIP), these results also lead to the conclusion that this dependence on coherence is qualitatively true and not simply a proof artifact.

There are several directions that can be explored in the future. First, as we have discussed in the introduction, block diagonal matrices are useful for modeling distributed measurement systems. It may therefore be of interest to specialize our results to some particular distributed systems. Take for example MIMO radar where multiple independent transmitters and receivers are arbitrarily distributed over an area of interest to sense targets in a scene. Data from the receivers with potentially high data rates needs to be sent to a central processor and to be coherently processed to achieve a maximum processing gain. The block diagonal structure studied here is potentially useful to analyze the possibility of compressing the data at the individual receivers before sending it to the central processor. For another example, block diagonal structure has also been exploited in observability studies of dynamical systems [29]. Our understanding of this and similar problems [2, 8, 9] may be enhanced using the results in this paper.

Acknowledgements

We would like to thank Borhan Sanandaji for helpful comments on an early version of the manuscript. We also gratefully acknowledge Justin Romberg and Alejandro Weinstein for valuable discussions and insightful comments.

References

- [1] H. Arguello and G. Arce. Restricted isometry property in coded aperture compressive spectral imaging. In *IEEE Statistical Signal Processing Workshop (SSP)*, 2012.
- [2] S.M. Asif and A.S. Charles. Estimation and dynamic updating of time-varying signals with sparse variations. In *Proceedings of International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2011.
- [3] R. Berinde, A.C. Gilbert, P. Indyk, H. Karloff, and M.J. Strauss. Combining geometry and combinatorics: A unified approach to sparse signal recovery. In *46th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, 2008.

- [4] E.J. Candès. Compressive sampling. In *Proceedings of International Congress of Mathematics*, 2006.
- [5] E.J. Candès. The restricted isometry property and its implications for compressed sensing. *Comptes Rendus Mathématique*, 346(9-10):589–592, 2008.
- [6] E.J. Candès, Y.C. Eldar, D. Needell, and P. Randall. Compressed sensing with coherent and redundant dictionaries. *Applied and Computational Harmonic Analysis*, 31(1):59–73, 2011.
- [7] E.J. Candès and M.B. Wakin. An introduction to compressive sampling. *IEEE Signal Processing Magazine*, 25(2):21–30, 2008.
- [8] A.S. Charles, S.M. Asif, J.K. Romberg, and C.J. Rozell. Sparsity penalties in dynamical system estimation. In *Proceedings of IEEE Conference on Information Sciences and Systems (CISS)*, 2011.
- [9] A.S. Charles and C.J. Rozell. Re-weighted ℓ_1 dynamic filtering for time-varying sparse signal estimation. *ArXiv preprint 1208.0325*, August 2012.
- [10] M.A. Davenport, P.T. Boufounos, M.B. Wakin, and R.G. Baraniuk. Signal processing with compressive measurements. *IEEE Journal of Selected Topics in Signal Processing*, 4(2):445–460, 2010.
- [11] M.A. Davenport and M.B. Wakin. Analysis of orthogonal matching pursuit using the restricted isometry property. *IEEE Transactions on Information Theory*, 56(9):4395–4401, 2010.
- [12] R.A. DeVore. Deterministic constructions of compressed sensing matrices. *Journal of Complexity*, 23(4-6):918–925, 2007.
- [13] D. Donoho and J. Tanner. Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 367(1906):4273–4293, 2009.
- [14] J. Haupt, W.U. Bajwa, G. Raz, and R. Nowak. Toeplitz compressed sensing matrices with applications to sparse channel estimation. *IEEE Transactions on Information Theory*, 56(11):5862–5875, 2010.
- [15] E.R. Kandel, J.H. Schwartz, and T.M. Jessell. *Principles of neural science*, volume 4. McGraw-Hill New York, 2000.
- [16] F. Krahmer, S. Mendelson, and H. Rauhut. Suprema of chaos processes and the restricted isometry property. *Arxiv preprint arXiv:1207.0235*, July 2012.
- [17] F. Krahmer and R. Ward. New and improved Johnson-Lindenstrauss embeddings via the restricted isometry property. *Arxiv preprint arXiv:1009.0744*, 2010.
- [18] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: isoperimetry and processes*, volume 23. Springer, 1991.

- [19] D. Needell and J.A. Tropp. CoSaMP: Iterative signal recovery from incomplete and inaccurate samples. *Applied and Computational Harmonic Analysis*, 26(3):301–321, 2009.
- [20] J.Y. Park, H.L. Yap, C.J. Rozell, and M.B. Wakin. Concentration of measure for block diagonal matrices with applications to compressive signal processing. *IEEE Transactions on Signal Processing*, 59(12):5859 –5875, 2011.
- [21] H. Rauhut. Compressive sensing and structured random matrices. In M. Fornasier, editor, *Theoretical Foundations and Numerical Methods for Sparse Recovery*, volume 9 of *Radon Series Comp. Appl. Math.*, pages 1–92. deGruyter, 2010.
- [22] H. Rauhut, J. Romberg, and J.A. Tropp. Restricted isometries for partial random circulant matrices. *Applied and Computational Harmonic Analysis*, 32(2):242 – 254, 2012.
- [23] C.J. Rozell, H.L. Yap, J.Y. Park, and M.B. Wakin. Concentration of measure for block diagonal matrices with repeated blocks. In *Proceedings of IEEE Conference on Information Sciences and Systems (CISS)*, 2010.
- [24] M. Rudelson and R. Vershynin. On sparse reconstruction from Fourier and Gaussian measurements. *Communications on Pure and Applied Mathematics*, 61(8):1025–1045, 2008.
- [25] M. Talagrand. *The generic chaining*. Springer, 2005.
- [26] A.M. Tillmann and M.E. Pfetsch. The computational complexity of RIP, NSP, and related concepts in compressed sensing. *Arxiv preprint arXiv:1205.2081*, 2012.
- [27] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. In Y.C. Eldar and G. Kutyniok, editors, *Compressed Sensing: Theory and Applications*. Cambridge University Press, 2012.
- [28] M.B. Wakin, J.Y. Park, H.L. Yap, and C.J. Rozell. Concentration of measure for block diagonal measurement matrices. In *Proceedings of International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2010.
- [29] M.B. Wakin, B.M. Sanandaji, and T.L. Vincent. On the observability of linear systems from random, compressive measurements. *IEEE Conf. Decision and Control (CDC)*, 2010.
- [30] J. Yang and Y. Zhang. Alternating direction algorithms for ℓ_1 -problems in compressive sensing. *SIAM Journal on Scientific Computing*, 33(1):250–278, 2011.
- [31] H.L. Yap, A. Eftekhari, M.B. Wakin, and C.J. Rozell. The restricted isometry property for block diagonal matrices. In *Proceedings of IEEE Conference on Information Sciences and Systems (CISS)*, 2011.
- [32] H.L. Yap and C.J. Rozell. Stable Takens’ embeddings for linear dynamical systems. *IEEE Transactions on Signal Processing*, 59(10):4781 –4794, 2011.
- [33] H.L. Yap, M.B. Wakin, and C.J. Rozell. Stable manifold embeddings with operators satisfying the restricted isometry property. In *Proceedings of IEEE Conference on Information Sciences and Systems (CISS)*, 2011.
- [34] Y. Zhang, J. Yang, and W. Yin. Yall1: Your algorithms for l1. yall1.blogs.rice.edu, 2011.

A Toolbox

This section collects a few general results that are used throughout the paper (mainly without proofs, for the benefit of the space).

Schatten norms possess the following useful property that mirrors Euclidean norms:

$$\text{Rank}(A)^{\frac{1}{p}-\frac{1}{q}} \|A\|_{S_q} \leq \|A\|_{S_p} \leq \|A\|_{S_q}, \quad (39)$$

for a matrix A , when $1 \leq q \leq p$. The following version of the Hölder inequality for matrices is used in this paper. For any pair of matrices A, B (such that AB exists), the following holds:⁶

$$\|AB\|_F \leq \|A\|_2 \|B\|_F. \quad (40)$$

For a random variable Z and $1 \leq p \leq q$, the following holds [21, Page 30]:

$$\mathbb{E}^p|Z| \leq \mathbb{E}^q|Z|. \quad (41)$$

Also, $\|C \cdot Z\|_{\psi_2} = |C| \|Z\|_{\psi_2}$ and, according to [27, Lemma 5.9], the following holds for a finite sequence of zero-mean independent random variables $\{Z_j\}$:

$$\left\| \sum_j Z_j \right\|_{\psi_2}^2 \lesssim \sum_j \|Z_j\|_{\psi_2}^2. \quad (42)$$

Throughout this section, let $g \in \mathbb{R}^N$ denote a vector whose entries are i.i.d. zero-mean unit variance sub-Gaussian random variables. For convenience, set $K := \max_{n \in [N]} \|g(n)\|_{\psi_2}$. For $t > 0$ and $n \in [N]$, the following holds by the definition of the sub-Gaussian norm [27]:

$$\log \mathbb{P}\{|g(n)| > t\} \lesssim -\frac{t^2}{\|g(n)\|_{\psi_2}^2}. \quad (43)$$

Also, the following holds when $N \gtrsim 1$, [27, Lemma 6.6]:

$$\sqrt{\mathbb{E} \|g\|_{\max}^2} = (\mathbb{E} \max_{n \in [N]} g^2(n))^{1/2} \lesssim K \sqrt{\log N}. \quad (44)$$

Furthermore, for $t \leq 1$, the following inequality is from [27, Corollary 5.17] and provides a lower bound on $\|g\|_2$:

$$\log \mathbb{P}\left\{\|g\|_2 < (1-t)\sqrt{N}\right\} \lesssim -\min\left(\frac{t^2}{K^4}, \frac{t}{K^2}\right)N. \quad (45)$$

Suppose that the entries of $G \in \mathbb{R}^{N \times J}$ are Gaussian random variables with zero-mean and unit variance. The next inequality provides an upper bound on the spectral norm of G , which directly follows from Corollary 5.35 in [27]. When $J \leq N$, and for $t > 0$, the following holds:

$$\log \mathbb{P}\left\{\|G\|_2 > (1+t)\sqrt{N} + \sqrt{J}\right\} \leq -t^2 N/2. \quad (46)$$

The next result essentially bounds the moments of a sum of independent random variables with those of a Rademacher sequence. The proof of this result (and the next one) uses the symmetrization

⁶ Let $\{b_j\}$ denote the columns of B . Then $\|AB\|_F^2 = \sum_j \|Ab_j\|_2^2 \leq \sum_j \|A\|_2^2 \|b_j\|_2^2 = \|A\|_2^2 \|B\|_F^2$.

technique, which has the following argument at its heart. If Z is a random variable taking values in a Banach space χ , then we can define its symmetrized version $Y = Z - Z'$, where Z' is an independent copy of Z . As suggested by the name, the distribution of Y is symmetric about the origin. In addition, the distributions of Y and ξY are the same, where ξ is a standard Bernoulli random variable that takes ± 1 with equal probability.

Lemma 8. [21, Lemma 6.7] *Let $\{Z_j\}$ be a finite sequence of independent random variables in the Banach space $(\chi, \|\cdot\|_\chi)$. Then, the following holds:*

$$\mathbb{E} \left\| \sum_j Z_j - \mathbb{E} Z_j \right\|_\chi \leq 2 \mathbb{E} \left\| \sum_j \xi_j Z_j \right\|_\chi.$$

The next lemma links the tail bounds of a sum of independent random variables and its symmetrized version. We remark that this result directly follows by applying equation 6.1 in [18] to $\sum_j Z_j - \mathbb{E} Z_j$.

Lemma 9. *Let $\{Z_j\}$ be defined as above, and let $\{Z'_j\}$ denote an independent copy of $\{Z_j\}$. The following holds for $t > 0$:*

$$\mathbb{P} \left\{ \left\| \sum_j Z_j - \mathbb{E} Z_j \right\|_\chi > 2 \mathbb{E} \left\| \sum_j Z_j - \mathbb{E} Z_j \right\|_\chi + t \right\} \leq 2 \mathbb{P} \left\{ \left\| \sum_j \xi_j (Z_j - Z'_j) \right\|_\chi > t \right\}.$$

In our proofs, we also require a (weak version of) the Khintchine inequality for operator norms that we state next.

Lemma 10. *If $\{A_l\}$, $l \in [L]$, is a sequence of matrices of the same dimension and rank of at most $K \gtrsim 1$, then the following holds.*

$$\mathbb{E} \left\| \sum_{l \in [L]} \xi_l A_l \right\|_2 \lesssim \sqrt{\log K} \left(\sum_{l \in [L]} \|A_l\|_2^2 \right)^{1/2}.$$

Proof. From [21, Theorem 6.14] and for every $2 \leq p < \infty$, we recall that

$$\mathbb{E}^p \left\| \sum_{l \in [L]} \xi_l A_l \right\|_{S_p} \lesssim \sqrt{p} \cdot \max \left(\left\| \left(\sum_{l \in [L]} A_l A_l^* \right)^{\frac{1}{2}} \right\|_{S_p}, \left\| \left(\sum_{l \in [L]} A_l^* A_l \right)^{\frac{1}{2}} \right\|_{S_p} \right).$$

The spectral norm is a special case of the Schatten norm with $p = \infty$. Therefore, the inequality above does not directly apply to our problem. As such, we need a more detailed argument here, which follows the approach of [27]. From (39), with $p = \infty$ and $q = \log J$, recall that

$$e^{-1} \|A\|_{S_{\log J}} \leq \|A\|_2 \leq \|A\|_{S_{\log J}}.$$

This equivalence, in combination with the fact that \mathbb{E}^p is increasing in p , i.e., (41), allows us to

write

$$\begin{aligned}
\mathbb{E} \left\| \sum_{l \in [L]} \xi_l A_l \right\|_2 &\leq \mathbb{E} \left\| \sum_{l \in [L]} \xi_l A_l \right\|_{S_{\log J}} \\
&\leq \mathbb{E}^{\log J} \left\| \sum_{l \in [L]} \xi_l A_l \right\|_{S_{\log J}} \\
&\lesssim \sqrt{\log J} \max \left(\left\| \left(\sum_{l \in [L]} A_l A_l^* \right)^{\frac{1}{2}} \right\|_{S_{\log J}}, \left\| \left(\sum_{l \in [L]} A_l^* A_l \right)^{\frac{1}{2}} \right\|_{S_{\log J}} \right) \\
&\leq e \sqrt{\log J} \max \left(\left\| \left(\sum_{l \in [L]} A_l A_l^* \right)^{\frac{1}{2}} \right\|_2, \left\| \left(\sum_{l \in [L]} A_l^* A_l \right)^{\frac{1}{2}} \right\|_2 \right) \\
&= e \sqrt{\log J} \max \left(\left\| \sum_{l \in [L]} A_l A_l^* \right\|_2^{\frac{1}{2}}, \left\| \sum_{l \in [L]} A_l^* A_l \right\|_2^{1/2} \right) \\
&\leq e \sqrt{\log J} \max \left(\left(\sum_{l \in [L]} \|A_l\|_2^2 \right)^{\frac{1}{2}}, \left(\sum_{l \in [L]} \|A_l\|_2^2 \right)^{\frac{1}{2}} \right) \\
&= e \sqrt{\log J} \cdot \left(\sum_{l \in [L]} \|A_l\|_2^2 \right)^{\frac{1}{2}},
\end{aligned}$$

as claimed. We assumed above that $J \geq e$ to produce the first line, and the second to the last line above uses the triangle inequality and the fact that $\|AA^*\|_2 = \|A\|_2^2$ for any matrix A . We remark that had we stopped at the fifth line above, we would have ended with the stronger original non-commutative Khintchine inequality for the spectral norm. However, the weaker bound given in the last line suffices for our purposes in this paper and completes the proof of Lemma 10. \square

We also list two trivial identities regarding covering numbers, which hold for every set \mathcal{S} , norm $\|\cdot\|$, and $r, a > 0$:

$$\begin{aligned}
\mathcal{N}(\mathcal{S}, a\|\cdot\|, r) &= \mathcal{N}(\mathcal{S}, \|\cdot\|, r/a) \\
\mathcal{N}(a\mathcal{S}, \|\cdot\|, r) &= \mathcal{N}(\mathcal{S}, \|\cdot\|, r/a).
\end{aligned} \tag{47}$$

For reference, we also provide an estimation for two integrals we encounter in the analysis:

$$\int_0^a \log \left(1 + \frac{b}{\nu} \right) d\nu \lesssim a \log \left(1 + \frac{b}{a} \right). \tag{48}$$

$$\int_0^a \sqrt{\log \left(1 + \frac{b}{\nu} \right)} d\nu \lesssim a \sqrt{\log \left(1 + \frac{b}{a} \right)}, \tag{49}$$

Both inequalities hold when $a \leq b$.

B Proof of Lemma 1

Equivalently, R can be constructed as follows. Let $\{r_{\tilde{n}}\}$, $\tilde{n} \in [\tilde{N}]$, denote the columns of R . The first column, r_1 , is chosen from the uniform distribution on the unit sphere in $\mathbb{R}^{\tilde{N}}$. For every $\tilde{n} \in [\tilde{N}] \setminus \{1\}$, $r_{\tilde{n}}$ is chosen from the uniform distribution on the unit sphere in the orthogonal complement of the span of the first $\tilde{n} - 1$ columns.

Let $g \in \mathbb{R}^{\tilde{N}}$ denote a standard Gaussian vector, that is a vector whose entries are i.i.d. zero-mean Gaussian random variables with unit variance. Since r_1 is drawn from the uniform distribution on the unit sphere in $\mathbb{R}^{\tilde{N}}$, the entries of r_1 have the same distribution as those in $g/\|g\|_2$. Since the distribution of R remains unchanged under permutation of its rows, every column of R has the same (marginal) distribution as $g/\|g\|_2$. This, alongside with the union bound, allows us to write the following for any $t > 0$:

$$\begin{aligned}
\mathbb{P}\left\{\mu(R) > t\sqrt{\log \tilde{N}}\right\} &= \mathbb{P}\left\{\max_{\tilde{n} \in [\tilde{N}]} \|r_{\tilde{n}}\|_{\max} > t\sqrt{\log \tilde{N}}/\sqrt{\tilde{N}}\right\} \\
&\leq \tilde{N} \cdot \max_{\tilde{n} \in [\tilde{N}]} \mathbb{P}\left\{\|r_{\tilde{n}}\|_{\max} > t\sqrt{\log \tilde{N}}/\sqrt{\tilde{N}}\right\} \\
&= \tilde{N} \cdot \mathbb{P}\left\{\max_{\tilde{n} \in [\tilde{N}]} \frac{|g(\tilde{n})|}{\|g\|_2} > t\sqrt{\frac{\log \tilde{N}}{\tilde{N}}}\right\} \\
&\leq \tilde{N}^2 \cdot \mathbb{P}\left\{\frac{|g(1)|}{\|g\|_2} > t\sqrt{\frac{\log \tilde{N}}{\tilde{N}}}\right\} \\
&= \tilde{N}^2 \cdot \mathbb{P}\left\{\frac{|g(1)|}{\|g\|_2} > \frac{t/2 \cdot \sqrt{\log \tilde{N}}}{(1 - 1/2)\sqrt{\tilde{N}}}\right\} \\
&\leq \tilde{N}^2 \cdot \mathbb{P}\left\{|g(1)| > t/2 \cdot \sqrt{\log \tilde{N}}\right\} + \tilde{N}^2 \cdot \mathbb{P}\left\{\|g\|_2 < (1 - 1/2)\sqrt{\tilde{N}}\right\} \\
&\lesssim \tilde{N}^2 e^{-\frac{C_2}{4\|g(1)\|_{\psi_2}^2} t^2 \log \tilde{N}} + \tilde{N}^2 e^{-C_3 \min\left(\frac{1}{4\|g(1)\|_{\psi_2}^4}, \frac{1}{2\|g(1)\|_{\psi_2}^2}\right) \tilde{N}} \\
&\leq \tilde{N}^2 e^{-\frac{C_2}{4\|g(1)\|_{\psi_2}^2} t^2 \log \tilde{N}} + \tilde{N}^2 e^{-\frac{C_3}{4\|g(1)\|_{\psi_2}^4} \tilde{N}},
\end{aligned}$$

where we used (43) and (45) to bound the failure probability. The last line holds because $\|g(1)\|_{\psi_2} = \sqrt{2/\pi}$.⁷ If we take $\tilde{N} \geq t^2 \log \tilde{N}$, we obtain that

$$\mathbb{P}\left\{\mu(R) > t\sqrt{\log \tilde{N}}\right\} \lesssim \tilde{N}^{2 - \frac{C_2}{4\|g(1)\|_{\psi_2}^2} t^2} + \tilde{N}^{2 - \frac{C_3}{4\|g(1)\|_{\psi_2}^4} t^2}. \quad (50)$$

We arrive at the advocated result when $t \gtrsim 1$:

$$\mathbb{P}\left\{\mu(R) > t\sqrt{\log \tilde{N}}\right\} \lesssim \tilde{N}^{-t} + \tilde{N}^{-t} = 2\tilde{N}^{-t}. \quad (51)$$

⁷ This is easily verified using the moments of the Gaussian distribution.

C Proof of Lemma 2

We use here the construction of R laid down in the beginning of B. As pointed out in the proof of Lemma 1, the columns of R are dependent but identically distributed as $g/\|g\|_2$, where g was defined there. Now, for every t , we can write

$$\begin{aligned} \mathbb{P}\left\{\gamma(R) \gtrsim 1 + \sqrt{\frac{J}{N}} + t\right\} &= \mathbb{P}\left\{\max_{\tilde{n} \in [\tilde{N}]} \|X(R, e_{\tilde{n}})\|_2 \gtrsim \frac{1}{\sqrt{N}} + \frac{1+t}{\sqrt{J}}\right\} \\ &\leq \tilde{N} \cdot \max_{\tilde{n} \in [\tilde{N}]} \mathbb{P}\left\{\|X(R, e_{\tilde{n}})\|_2 \gtrsim \frac{1}{\sqrt{N}} + \frac{1+t}{\sqrt{J}}\right\} \\ &= \tilde{N} \cdot \mathbb{P}\left\{\|X(R, e_1)\|_2 \gtrsim \frac{1}{\sqrt{N}} + \frac{1+t}{\sqrt{J}}\right\}. \end{aligned} \quad (52)$$

The second line uses the union bound and the last line holds due to the identical distribution of the columns of R . It remains to find an upper bound for the probability in the last line above. Recall that r_1 has the same distribution as $g/\|g\|_2$, and thus $X(R, e_1)$ has the same distribution as $G/\|G\|_F$, where $G \in \mathbb{C}^{N \times J}$ is formed by reshaping g . Therefore, $\|X(R, e_1)\|_2$ has the same distribution as $\|G\|_2 / \|G\|_F$. For fixed $t \leq 1$, the following convenient inequality holds:⁸

$$\frac{1}{\sqrt{N}} + \frac{1+t}{\sqrt{J}} \gtrsim \frac{(1+t/3)\sqrt{N} + \sqrt{J}}{(1-t/3)\sqrt{\tilde{N}}}. \quad (53)$$

Now, we can write that

$$\begin{aligned} &\mathbb{P}\left\{\|X(R, e_1)\|_2 \gtrsim \frac{1}{\sqrt{N}} + \frac{1+t}{\sqrt{J}}\right\} \\ &= \mathbb{P}\left\{\frac{\|G\|_2}{\|G\|_F} \gtrsim \frac{1}{\sqrt{N}} + \frac{1+t}{\sqrt{J}}\right\} \\ &\leq \mathbb{P}\left\{\frac{\|G\|_2}{\|G\|_F} \gtrsim \frac{(1+\frac{t}{3})\sqrt{N} + \sqrt{J}}{(1-\frac{t}{3})\sqrt{\tilde{N}}}\right\} \\ &\leq \mathbb{P}\left\{\|G\|_2 \gtrsim \left(1 + \frac{t}{3}\right)\sqrt{N} + \sqrt{J}\right\} + \mathbb{P}\left\{\|G\|_F \lesssim \left(1 - \frac{t}{3}\right)\sqrt{\tilde{N}}\right\} \\ &\lesssim e^{-\frac{1}{18}t^2N} + e^{-C_3 \min\left(\frac{t^2}{9\|g(1)\|_{\psi_2}^4}, \frac{t}{3\|g(1)\|_{\psi_2}^2}\right)\tilde{N}} \\ &\leq e^{-\frac{1}{18}t^2N} + e^{-\frac{C_3}{9\|g(1)\|_{\psi_2}^4}t^2\tilde{N}} \\ &\leq 2e^{-C_4t^2N}, \end{aligned}$$

where the second line uses (53). The fourth line uses the inequalities (45) and (46) and the fact that $\|G\|_F = \|g\|_2$. The second to last line follows because $\|g(1)\|_{\psi_2} = \sqrt{2/\pi}$ and $t \leq 1$. The above upper bound in combination with (52) leads us to the following conclusion:

$$\mathbb{P}\left\{\gamma(R) \gtrsim 1 + \sqrt{\frac{J}{N}} + t\right\} \lesssim \tilde{N}e^{-C_4t^2N}.$$

⁸ This inequality easily follows from the fact that $1+t \geq (1+t/3)(1-t/3)^{-1}$ holds for every $t \leq 1$.

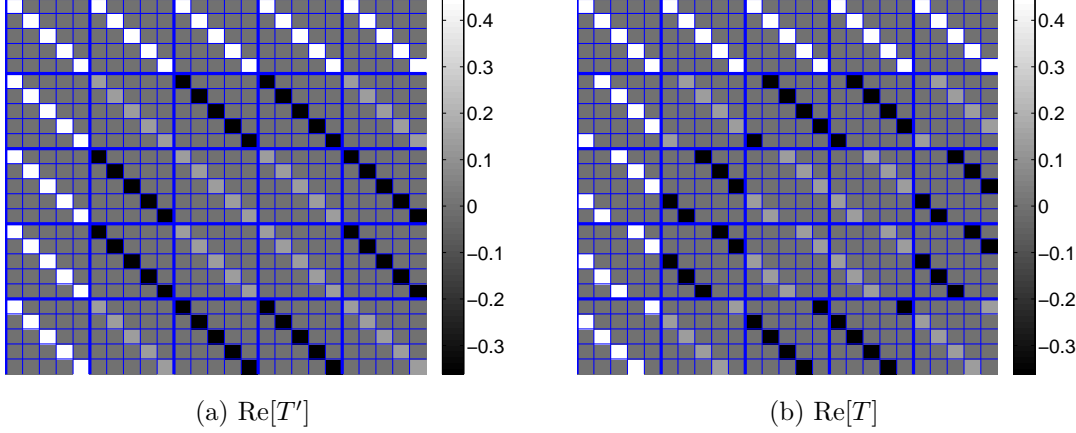


Figure 3: A visual illustration of the transformation from T' to T for $P = J = 5$. See the explanation in D.

We complete the proof of Lemma 2 by taking $N \geq 2C_4^{-1}t^{-2} \log \tilde{N}$.

D Proof of Lemma 3

Let $T' := F_J \otimes I_P \in \mathbb{C}^{PJ \times PJ}$, where \otimes stands for the Kronecker product. Here, F_J and I_P are the Fourier and canonical orthobases of size, respectively, $J \times J$ and $P \times P$. We consider the natural partitioning of T' into J row-submatrices of size $P \times PJ$, denoted by T'_j , $j \in [J]$. Now, for every $j \in [J]$, cyclically shift the rows of T'_j by $j - 1$ times upward and create $T_j \in \mathbb{C}^{P \times PJ}$. Then form the matrix $T \in \mathbb{C}^{PJ \times PJ}$ by replacing every T'_j with T_j . The transformation of T' to T is visualized in Figure 3.

Using the properties of the Kronecker product, it is easily verified that T' is an orthobasis for \mathbb{C}^{PJ} . Due to its structure, we also observe that the nonzero entries of the p_1 th and p_2 th columns of T do not overlap when $(p_1 - p_2) \bmod P \neq 0$, and so these columns are orthogonal. Otherwise, when $(p_1 - p_2) \bmod P = 0$, we need a more subtle argument. Under this condition, the inner product of the p_1 th and p_2 th columns of T' equals the inner product of the $\lceil p_1/P \rceil$ th and $\lceil p_2/P \rceil$ th columns of F_J and is indeed zero. The inner product of the p_1 th and p_2 th columns remains zero under the transformation of T' to T since this transformation only amounts to a permutation in the rows of T' . Therefore, T is an orthobasis for \mathbb{C}^{PJ} . As for computing $\gamma(T)$, the structure of T guarantees that, for each j , every column and row of $X(e_j, T)$ has only one nonzero entry with the magnitude of $1/\sqrt{J}$. Therefore, $\gamma(T) = 1$.

Finally, it can be easily verified that $\hat{x}/\sqrt{J} = Tx_e$, where $x_e = [x^T, 0, 0, \dots, 0]^T \in \mathbb{C}^{PJ}$, and by (16), $\Gamma x = \Xi Tx_e$. If x is S -sparse, so is x_e . This completes the proof of Lemma 3.

E Proof of Lemma 6

The arguments used in this section are largely adapted from [24]. In what follows, we let $\mathcal{B}_A^{\tilde{N}}$ denote the unit ball with respect to the norm $\|\cdot\|_A$ in $\mathbb{C}^{\tilde{N}}$. Also, $\mathcal{B}_1^{\tilde{N}}$ and $\mathcal{B}_2^{\tilde{N}}$, respectively, denote the unit ℓ_1 -ball and ℓ_2 -ball in $\mathbb{C}^{\tilde{N}}$. For $T \subseteq [\tilde{N}]$, we let \mathcal{B}_A^T denote the unit ball in the $\#T$ -dimensional

subspace of $\mathbb{C}^{\tilde{N}}$ spanned by $\{e_{\tilde{n}}\}$, $\tilde{n} \in T$. We define \mathcal{B}_1^T and \mathcal{B}_2^T similarly.

The first thing to notice is that when $\alpha \in \Omega_S/\sqrt{S}$, then $\|\alpha\|_1 \leq 1$. From the hypothesis of the lemma, we then have that

$$\|\alpha\|_A \leq \frac{\kappa}{\sqrt{\widetilde{M}}}. \quad (54)$$

This implies that for every support $T \subset [\tilde{N}]$ with $\#T = S$, we have

$$\frac{\mathcal{B}_2^T}{\sqrt{S}} \subseteq \frac{\kappa}{\sqrt{\widetilde{M}}} \cdot \mathcal{B}_A^T. \quad (55)$$

On the other hand, Ω_S/\sqrt{S} can be equivalently represented as

$$\frac{\Omega_S}{\sqrt{S}} = \bigcup_{\#T=S} \frac{\mathcal{B}_2^T}{\sqrt{S}}. \quad (56)$$

Together, (55) and (56) imply that

$$\frac{\Omega_S}{\sqrt{S}} \subseteq \bigcup_{\#T=S} \frac{\kappa}{\sqrt{\widetilde{M}}} \cdot \mathcal{B}_A^T. \quad (57)$$

We also record that

$$\frac{\Omega_S}{\sqrt{S}} \subseteq \frac{\kappa}{\sqrt{\widetilde{M}}} \cdot \mathcal{B}_{X^*}^{\tilde{N}}, \quad (58)$$

which dictates that

$$\mathcal{N}(\Omega_S/\sqrt{S}, \|\cdot\|_A, \nu) = 1 \quad (59)$$

if $\nu \geq \kappa/\sqrt{\widetilde{M}}$. This proves the second statement in Lemma 6. Otherwise, if $\nu < \kappa/\sqrt{\widetilde{M}}$, we continue with the rest of the argument. In light of (57), we have that

$$\begin{aligned} \mathcal{N}\left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_A, \nu\right) &\leq \mathcal{N}\left(\bigcup_{\#T=S} \frac{\kappa}{\sqrt{\widetilde{M}}} \cdot \mathcal{B}_A^T, \|\cdot\|_A, \nu\right) \\ &\leq \binom{\tilde{N}}{S} \cdot \mathcal{N}\left(\frac{\kappa}{\sqrt{\widetilde{M}}} \cdot \mathcal{B}_A^T, \|\cdot\|_A, \nu\right) \\ &= \binom{\tilde{N}}{S} \cdot \mathcal{N}\left(\mathcal{B}_A^T, \|\cdot\|_A, \nu\kappa^{-1}\sqrt{\widetilde{M}}\right), \end{aligned} \quad (60)$$

where the last line uses the second inequality in (47). An estimate for the covering number in the last line above is available for $r > 0$, namely⁹

$$\mathcal{N}(\mathcal{B}_A^T, \|\cdot\|_A, r) \leq (1 + 2/r)^{2\#T}. \quad (61)$$

⁹ This is proved similar to Lemma 5.2 in [27], after exchanging the Euclidean metric with $\|\cdot\|_A$ and accounting for the complex vector space.

Plugging the bound above into (60), we arrive at

$$\mathcal{N}\left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_A, \nu\right) \leq \binom{\tilde{N}}{S} \cdot \left(1 + \frac{2\kappa}{\nu\sqrt{\widetilde{M}}}\right)^{2S} \leq \left(\frac{e\tilde{N}}{S}\right)^S \cdot \left(1 + \frac{2\kappa}{\nu\sqrt{\widetilde{M}}}\right)^{2S}. \quad (62)$$

The last inequality holds because $\binom{n}{m} \leq (en/m)^m$ for any pair of integers $m \leq n$. When $\tilde{N} \gtrsim 1$, (62) implies that

$$\begin{aligned} \log \mathcal{N}\left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_A, \nu\right) &\leq S \log\left(\frac{e\tilde{N}}{S}\right) + 2S \log\left(1 + \frac{2\kappa}{\nu\sqrt{\widetilde{M}}}\right) \\ &\lesssim S \log \tilde{N} + S \log\left(1 + \frac{2\kappa}{\nu\sqrt{\widetilde{M}}}\right). \end{aligned} \quad (63)$$

The bound above is less effective for small values of ν . To seek a second bound on the covering number, we begin from the containment

$$\frac{\Omega_S}{\sqrt{S}} \subseteq \mathcal{B}_1^{\tilde{N}}, \quad (64)$$

which immediately dictates that

$$\mathcal{N}\left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_A, \nu\right) \leq \mathcal{N}\left(\mathcal{B}_1^{\tilde{N}}, \|\cdot\|_A, \nu\right). \quad (65)$$

In order to compute the covering number on the right hand side of (65), we use the following result, which is proved in F.

Lemma 11. *Let $B_{1,\tilde{N}}$ denote the ℓ_1 -ball in $\mathbb{R}^{\tilde{N}}$, and consider the norm $\|\cdot\|_A$ on $\mathbb{C}^{\tilde{N}}$, which satisfies the hypothesis of Lemma 6. Naturally, $\|\cdot\|_A$ induces a norm on $\mathbb{R}^{\tilde{N}} \subset \mathbb{C}^{\tilde{N}}$ (which is represented with the same notation for convenience). For $\nu > 0$ and $\widetilde{M} \gtrsim 1$, it holds that*

$$\log\left(\mathcal{N}\left(B_{1,\tilde{N}}, \|\cdot\|_A, \nu\right)\right) \lesssim \frac{\kappa^2}{\nu^2 \widetilde{M}} \cdot \log^2 \tilde{N}.$$

Now consider an arbitrary $\beta \in \mathcal{B}_1^{\tilde{N}}$, and note that $\text{Re}[\beta], \text{Im}[\beta] \in B_{1,\tilde{N}}$. Let

$$\mathcal{C}_1 := \mathcal{C}\left(B_{1,\tilde{N}}, \|\cdot\|_A, \nu/2\right)$$

denote a minimal $(\nu/2)$ -cover for $B_{1,\tilde{N}}$ with respect to the metric $\|\cdot\|_A$. Therefore, there exist $p_1, p_2 \in \mathcal{C}_1$ such that $\|\text{Re}[\beta] - p_1\|_A, \|\text{Im}[\beta] - p_2\|_A \leq \nu/2$. It follows by the triangle inequality that

$$\begin{aligned} \|\beta - (p_1 + ip_2)\|_A &= \|(\text{Re}[\beta] - p_1) + i(\text{Im}[\beta] - p_2)\|_A \\ &\leq \|\text{Re}[\beta] - p_1\|_A + \|\text{Im}[\beta] - p_2\|_A \\ &\leq \nu. \end{aligned}$$

Therefore, $\{p_1 + ip_2 : p_1, p_2 \in \mathcal{C}_1\}$ is a cover for $\mathcal{B}_1^{\tilde{N}}$, and clearly

$$\mathcal{N}(\mathcal{B}_1^{\tilde{N}}, \|\cdot\|_A, \nu) \leq \left(\mathcal{N}(B_{1,\tilde{N}}, \|\cdot\|_A, \nu/2) \right)^2.$$

It now follows from (65), Lemma 11, and the inequality above that

$$\log \mathcal{N} \left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_A, \nu \right) \lesssim \frac{\kappa^2}{\nu^2 \widetilde{M}} \cdot \log^2 \tilde{N}. \quad (66)$$

Combining (63) and (66), we finally arrive at

$$\log \mathcal{N} \left(\frac{\Omega_S}{\sqrt{S}}, \|\cdot\|_A, \nu \right) \lesssim \min \left(S \log \tilde{N} + S \log \left(1 + \frac{2\kappa}{\nu \sqrt{\widetilde{M}}} \right), \frac{\kappa^2}{\nu^2 \widetilde{M}} \cdot \log^2 \tilde{N} \right),$$

which holds when $\widetilde{M} \gtrsim 1$. This completes the proof of Lemma 6.

F Proof of Lemma 11

Consider an arbitrary $\beta \in B_{1,\tilde{N}}$. Also consider a random vector Z that takes a value in $\{0\} \cup \{\text{sgn}(\beta(\tilde{n})) \cdot e_{\tilde{n}}\}$, $\tilde{n} \in [\tilde{N}]$. It takes $\text{sgn}(\beta(\tilde{n})) \cdot e_{\tilde{n}}$ with probability of $|\beta(\tilde{n})|$ and zero otherwise. Clearly, $\mathbb{E}Z = \beta$. Now, we wish to approximate β with the average of L independent copies of Z , denoted by $\{Z_l\}$, $l \in [L]$. The expected approximation error, measured in $\|\cdot\|_A$, would be

$$\mathbb{E} \left\| \beta - \frac{1}{L} \sum_{l \in [L]} Z_l \right\|_A.$$

Since the argument of the norm is zero-mean, we can use the symmetrization technique by invoking Lemma 8 from the Toolbox to obtain that

$$\mathbb{E} \left\| \beta - \frac{1}{L} \sum_{l \in [L]} Z_l \right\|_A = \frac{1}{L} \mathbb{E} \left\| \sum_{l \in [L]} Z_l - \mathbb{E} Z_l \right\|_A \leq \frac{2}{L} \mathbb{E} \left\| \sum_{l \in [L]} \xi_l Z_l \right\|_A, \quad (67)$$

where $\{\xi_l\}$, $l \in [L]$, is a Rademacher sequence. According to Lemma 6, $\|\cdot\|_A = \|A(\cdot)\|_2$ and $A(\cdot)$ is a linear map. Therefore

$$\mathbb{E} \left\| \sum_{l \in [L]} \xi_l Z_l \right\|_A = \mathbb{E} \left\| \sum_{l \in [L]} \xi_l A(Z_l) \right\|_2.$$

Also, from the hypothesis of Lemma 6, $\text{Rank}(A(Z_l)) \leq \widetilde{M}$ for every l . We now invoke a Khintchine-type inequality for the operator norm (stated in Lemma 10 from the Toolbox), which allows us to

continue our argument as

$$\begin{aligned}
\mathbb{E} \left\| \sum_{l \in [L]} \xi_l A(Z_l) \right\|_2 &= \mathbb{E}^Z \mathbb{E}^\xi \left\| \sum_{l \in [L]} \xi_l A(Z_l) \right\|_2 \\
&\lesssim \sqrt{\log \widetilde{M}} \cdot \mathbb{E}^Z \left(\sum_{l \in [L]} \|A(Z_l)\|_2^2 \right)^{1/2} \\
&\leq \sqrt{\log \widetilde{M}} \cdot \left(\sum_{l \in [L]} \mathbb{E} \|A(Z_l)\|_2^2 \right)^{1/2} \\
&= \sqrt{\log \widetilde{M}} \cdot \left(\sum_{l \in [L]} \sum_{\tilde{n} \in [\tilde{N}]} |\beta_{\tilde{n}}| \cdot \|A(e_{\tilde{n}})\|_2^2 \right)^{1/2} \\
&\leq \sqrt{\log \widetilde{M}} \cdot \left(L \cdot \max_{\tilde{n} \in \tilde{N}} \|A(e_{\tilde{n}})\|^2 \right)^{1/2} \\
&= \sqrt{L \log \widetilde{M}} \cdot \max_{\tilde{n} \in \tilde{N}} \|e_{\tilde{n}}\|_A \\
&\leq \kappa \sqrt{\frac{L}{\widetilde{M}} \log \widetilde{M}}, \tag{68}
\end{aligned}$$

where, conditioned on $\{Z_l\}$, $l \in [L]$, the Khintchine inequality was used to produce the second line. The third line follows from the Jensen inequality. The fifth line uses the assumption that $\beta \in B_{1, \tilde{N}}$. The last line follows from the hypothesis of Lemma 6. Using the above inequality in combination with (67) yields

$$\mathbb{E} \left\| \beta - \frac{1}{L} \sum_{l \in [L]} Z_l \right\|_X \lesssim \kappa \sqrt{\frac{\log \widetilde{M}}{L \widetilde{M}}}.$$

To keep the average no larger than ν , it suffices to take

$$L \gtrsim \frac{\kappa^2}{\nu^2 \widetilde{M}} \cdot \log \widetilde{M}.$$

With this choice of L , there exists a linear combination of independent copies of Z that falls within a ν distance of β . In other words, we have shown that for an arbitrary $\beta \in B_{1, \tilde{N}}$, there exists an average of L elements of $\{0\} \cup \{\pm e_{\tilde{n}}\}$ that is of distance at most ν from β . There are $2\tilde{N} + 1$ elements in the aforementioned set, and so there are $(2\tilde{N} + 1)^L$ possibilities for the average. We therefore conclude that

$$\mathcal{N} \left(B_{1, \tilde{N}}, \|\cdot\|_A, \nu \right) \leq \left(2\tilde{N} + 1 \right)^{C_5 \log \widetilde{M} \cdot \kappa^2 / \nu^2 \widetilde{M}},$$

or

$$\begin{aligned}
\log \mathcal{N} \left(B_{1, \tilde{N}}, \|\cdot\|_A, \nu \right) &\lesssim \frac{\kappa^2}{\nu^2 \widetilde{M}} \cdot \log \widetilde{M} \log (2\tilde{N} + 1) \\
&\lesssim \frac{\kappa^2}{\nu^2 \widetilde{M}} \cdot \log \widetilde{M} \log \tilde{N} \\
&\leq \frac{\kappa^2}{\nu^2 \widetilde{M}} \cdot \log^2 \tilde{N}
\end{aligned}$$

when $\tilde{N} \gtrsim 1$. This completes the proof of Lemma 11.